

Asynchronous Adaptive Networks

Ali H. Sayed and Xiaochuan Zhao

Electrical Engineering Department
University of California, Los Angeles

A. H. Sayed and X. Zhao, "Asynchronous adaptive networks," in *Cognitive Dynamic Systems*, S. Haykin, Ed., Wiley, NY, 2016.

Abstract—In a recent article [1] we surveyed advances related to adaptation, learning, and optimization over synchronous networks. Various distributed strategies were discussed that enable a collection of networked agents to interact locally in response to streaming data and to continually learn and adapt to track drifts in the data and models. Under reasonable technical conditions on the data, the adaptive networks were shown to be mean-square stable in the slow adaptation regime, and their mean-square-error performance and convergence rate were characterized in terms of the network topology and data statistical moments [2]. Classical results for single-agent adaptation and learning were recovered as special cases. Following the works [3]–[5], this chapter complements the exposition from [1] and extends the results to asynchronous networks. The operation of this class of networks can be subject to various sources of uncertainties that influence their dynamic behavior, including randomly changing topologies, random link failures, random data arrival times, and agents turning on and off randomly. In an asynchronous environment, agents may stop updating their solutions or may stop sending or receiving information in a random manner and without coordination with other agents. The presentation will reveal that the mean-square-error performance of asynchronous networks remains largely unaltered compared to synchronous networks. The results justify the remarkable resilience of cooperative networks in the face of random events.

I. INTRODUCTION: COGNITION AND LEARNING

Nature is laden with examples where complex and sophisticated patterns of behavior emanate from limited interactions among simple elements, and from the aggregation and processing of decentralized pieces of information collected by dispersed agents over a graph. Examples abound in the realm of biological networks, where remarkable patterns of coordinated behavior manifest themselves, for example, in the form of fish schooling, bird formations, or bee swarming [6]. While each individual agent in these networks is incapable of complex decision-making on its own, it is the continuous coordination and sharing of information among neighboring agents that lead to effective multi-agent formations.

Network science is the field that deals with issues related to the aggregation, processing, and diffusion of information over graphs linking a multitude of agents. While the interactions over such graphs can be studied and characterized from the

perspective of cluster formations, degrees of connectivity, and small-world effects [7]–[9], it is the possibility of having agents interact dynamically with each other, and influence each other's behavior, that opens up a plethora of notable questions. For example, a better understanding of how local interactions influence the global pattern of behavior at the network level can lead to a broader understanding of how localized interactions in the social sciences, life sciences, and system sciences influence the evolution of the respective networks. For long, system theory has focused on studying stand-alone dynamic systems with great success. However, rapid advances in the biological sciences, animal behavior studies, and in the neuroscience of the brain, are revealing the striking power of coordination among networked units (e.g., [6], [10]–[12]). These discoveries are motivating greater efforts towards a deeper examination of information processing over graphs in several disciplines including signal processing, machine learning, optimization, and control (see, e.g., [2], [13] and the references therein).

The earlier article [1] surveys the field of synchronous adaptive networks and how collaboration among agents can lead to superior adaptation and learning performance over graphs. The monograph [2] provides a deeper treatment of the subject matter, including derivations and various additional aspects of network behavior. Adaptive networks consist of a collection of agents with learning abilities. The agents interact with each other on a local level and diffuse information across the network to solve inference and optimization tasks in a decentralized manner. Such networks are scalable, robust to node and link failures, and are particularly suitable for learning from big data sets by tapping into the power of collaboration among distributed agents. The networks are also endowed with cognitive abilities due to the sensing abilities of their agents, their interactions with their neighbors, and the embedded feedback mechanisms for acquiring and refining information. Each agent is not only capable of sensing data and experiencing the environment directly, but it also receives information through interactions with its neighbors and processes and analyzes this information to drive its learning process.

As already indicated in [1], [2], there are many good reasons for the peaked interest in networked solutions, especially in this day and age when the word "network" has become commonplace whether one is referring to social networks, power networks, transportation networks, biological networks, or other networks. Some of these reasons have to do with the benefits of cooperation over networks in terms of improved

This article is scheduled to appear as a book chapter in the edited volume *Cognitive Dynamic Systems*, S. Haykin, Ed., Wiley, NY, 2016. The authors are with the Electrical Engineering Department, University of California, Los Angeles. Emails: sayed@ucla.edu and xiaochuanzhao@ucla.edu. This work was supported in part by NSF grants CCF-1011918 and ECCS-1407712. The authors are grateful to IEEE for allowing reproduction of the material from [1] in this book chapter.

performance and improved robustness and resilience to failure. Other reasons deal with privacy and secrecy considerations where agents may not be comfortable sharing their data with remote fusion centers. In other situations, the data may already be available in dispersed locations, as happens with cloud computing. One may also be interested in learning and extracting information through data mining from large data sets. Decentralized learning procedures offer an attractive approach to dealing with such data sets. Decentralized mechanisms can also serve as important enablers for the design of robotic swarms, which can assist in the exploration of disaster areas.

Motivated by these observations, we devote reasonable effort towards clarifying the limits of performance of networked solutions by relying on statistical analysis tools. Whenever necessary, derivations are provided to complement the discussion with references to the pertinent literature for the longer arguments that are omitted for space considerations. The results are illustrated by means of examples dealing with applications involving distributed estimation, learning, optimization, and adaptation. For other applications in the areas of intrusion detection, online dictionary learning, target localization, spectrum sensing, sparse data recovery, and biological networks, the reader may refer to [2], [14]–[18] and the references therein.

A. Asynchronous Behavior

The presentation in [1] focused on the case of synchronous networks where data arrive at all agents in a synchronous manner and updates by the agents are also performed in a synchronous manner. The network topology was assumed to remain largely static during the adaptation process. Under these conditions, the limits of performance and stability of these networks were identified in some detail [1], [2] for two main classes of distributed strategies including consensus and diffusion constructions. In this chapter, we extend the overview from [1] to cover *asynchronous* environments as well. In such environments, the operation of the network can suffer from the interference of various random events including randomly changing topologies, random link failures, random data arrival times, and agents turning on and off randomly. Agents may stop updating their solutions or may stop sending or receiving information in a random manner and without coordination with other agents. Results in [3]–[5] examined the implications of asynchronous events on network performance in some detail and under a fairly general model for the random events. The purpose of this chapter is to summarize the key conclusions from these works in a manner that complements the presentation from [1] for the benefit of the reader. While the works [3]–[5] consider a broader formulation involving complex-valued arguments, we shall limit the discussion here to real-valued arguments in order not to overload the notation. Proofs and derivations are omitted and can be found in the above references; the emphasis is on presenting the results in a motivated manner and on commenting on the relevance of these results and the insights they provide into the operation of asynchronous networks.

We indicated in [3]–[5] that there already exist many useful studies in the literature on the performance of consensus strate-

gies in the presence of asynchronous events (see, e.g., [19]–[29]). There are also some studies in the context of diffusion strategies [30], [31]. However, with the exception of the latter two works, the earlier references assumed conditions that are not generally favorable for applications involving continuous *adaptation* and learning. For example, some of the works assumed a decaying step-size, which turns off adaptation after sufficient iterations have passed. Some other works assumed noise free data, which is a hindrance when learning from data perturbed by interferences and distortions. A third class of works focused on studying pure averaging algorithms, which are not required to respond to continuous data streaming. In the works [3]–[5], we adopted a more general asynchronous model that removes these limitations by allowing for various sources of random events and, moreover, the events are allowed to occur simultaneously. We also examined learning algorithms that respond to streaming data to enable adaptation. The main conclusion from the analysis in these works, and which will be summarized in future sections, is that asynchronous networks can still behave in a mean-square stable manner for sufficiently small step-sizes and, interestingly, their steady-state performance level is only slightly affected in comparison to synchronous behavior. The iterates computed by the various agents are still able to converge and hover around an agreement state with a small mean-square-error. These are reassuring results that support the intrinsic robustness and resilience of network-based cooperative solutions.

B. Organization

Readers will benefit more from this chapter if they review first the earlier article [1]. We continue to follow a similar structure here, as well as a similar notation, since the material in both this chapter and the earlier reference [1] are meant to complement each other.

In order to highlight the main features of adaptive networks, we organize the presentation of the article into three main components. The first part (Sec. II) reviews fundamental results on adaptation and learning by *single* stand-alone agents. The emphasis is on stochastic-gradient constructions, and also on asynchronous implementations. A general formulation is considered that allows us to extract classical results for synchronous adaptive filtering as special cases. The level of generality considered in this section is meant to bring forth commonalities that exist among several domains relating to adaptation, learning, and optimization.

The second part (Sec. III) of the chapter covers asynchronous centralized solutions. The objective is to explain the gain in performance that results from aggregating the data from the agents and processing it centrally at a fusion center. The centralized performance is used as a frame of reference for assessing various implementations. While centralized solutions can be powerful, they nevertheless suffer from a number of limitations. First, in real-time applications where agents collect data continuously, the repeated exchange of information back and forth between the agents and the fusion center can be costly especially when these exchanges occur over wireless links or require nontrivial routing resources. Second, in some

sensitive applications, agents may be reluctant to share their data with remote centers for various reasons including privacy and secrecy considerations. More importantly perhaps, centralized solutions have a critical point of failure: if the central processor fails, then this solution method collapses altogether.

For these reasons, we cover in the remaining sections of the chapter (Secs. IV and V) distributed asynchronous strategies of the consensus and diffusion types, and examine their dynamics, stability, and performance metrics. In the distributed mode of operation, agents are connected by a topology and they are permitted to share information only with their immediate neighbors. The study of the behavior of such networked agents is more challenging than in the single-agent and centralized modes of operation due to the coupling among interacting agents and due to the fact that the networks are generally sparsely connected. The presentation in the chapter clarifies the effect of network topology on performance and leads to results that enable the designer to compare various strategies against each other and against the centralized solution.

II. SINGLE-AGENT ADAPTATION AND LEARNING

We begin our treatment by reviewing stochastic gradient algorithms, with emphasis on their application to the problems of adaptation and learning by stand-alone agents. We will be using the term “learning” to refer broadly to the ability of an agent to extract information about some unknown parameter from streaming data, such as estimating the parameter itself or learning about some of its features. We will be using the term “adaptation” to refer broadly to the ability of the learning algorithm to track drifts in the parameter, which are usually reflected in changes in the statistical properties of the observed data. The two attributes of learning and adaptation will be embedded simultaneously into the algorithms discussed in this work. We will also be using the term “streaming data” regularly because we are interested in algorithms that perform continuous learning and adaptation and that, therefore, are able to improve their performance in response to continuous streams of data arriving at the agent(s). This is in contrast to off-line algorithms, where the data are first aggregated before being processed for extraction of information. The presentation in this section summarizes some classical results on stochastic-gradient algorithms for adaptation and learning, and provides some additional insights that are useful for our later study of the more demanding scenario of adaptation and learning by a collection of networked agents.

A. Risk and Loss Functions

Thus, let $J(w) : \mathbb{R}^{M \times 1} \mapsto \mathbb{R}$ denote a real-valued (cost or utility or risk) function of a real-valued vector argument, $w \in \mathbb{R}^{M \times 1}$. The variable w can be complex-valued, and many of the results in this work can be extended to the complex domain as well. However, some important technical differences arise when dealing with complex arguments. These differences are beyond the scope of this chapter and they are addressed in [2]–[4], along with other relevant topics. It is sufficient for our purposes here to convey the main ideas by limiting the presentation to real arguments without much loss in generality.

We denote the gradient vectors of $J(w)$ relative to w and w^\top by the following row and column vectors, respectively, where the first expression is also referred to as the Jacobian of $J(w)$ relative to w :

$$\nabla_w J(w) \triangleq \left[\frac{\partial J(w)}{\partial w_1}, \frac{\partial J(w)}{\partial w_2}, \dots, \frac{\partial J(w)}{\partial w_M} \right] \quad (1a)$$

$$\nabla_{w^\top} J(w) \triangleq [\nabla_w J(w)]^\top \quad (1b)$$

These definitions are in terms of the partial derivatives of $J(w)$ relative to the individual entries of $w = \text{col}\{w_1, w_2, \dots, w_M\}$, where the notation $\text{col}\{\cdot\}$ refers to a column vector that is formed by stacking its arguments on top of each other. Likewise, the Hessian matrix of $J(w)$ with respect to w is defined as the following $M \times M$ symmetric matrix:

$$\nabla_w^2 J(w) \triangleq \nabla_{w^\top} [\nabla_w J(w)] = \nabla_w [\nabla_{w^\top} J(w)] \quad (1c)$$

which is constructed from two successive gradient operations. It is common in adaptation and learning applications for the risk function $J(w)$ to be constructed as the expectation of some loss function, $Q(w; \mathbf{x})$, where the **boldface** variable \mathbf{x} is used to denote some random data, say,

$$J(w) = \mathbb{E} Q(w; \mathbf{x}) \quad (2)$$

and the expectation is evaluated over the distribution of \mathbf{x} .

Example II.1 (Mean-square-error costs). Let d denote a zero-mean scalar random variable with variance $\sigma_d^2 = \mathbb{E} d^2$, and let \mathbf{u} denote a zero-mean $1 \times M$ random vector with covariance matrix $R_u = \mathbb{E} \mathbf{u}^\top \mathbf{u} > 0$. The combined quantities $\{d, \mathbf{u}\}$ represent the random variable \mathbf{x} referred to in (2). The cross-covariance vector is denoted by $r_{du} = \mathbb{E} d \mathbf{u}^\top$. We formulate the problem of estimating d from \mathbf{u} in the linear least-mean-squares sense or, equivalently, the problem of seeking the vector w° that minimizes the quadratic cost function:

$$J(w) \triangleq \mathbb{E} (d - \mathbf{u}w)^2 = \sigma_d^2 - r_{du}^\top w - w^\top r_{du} + w^\top R_u w \quad (3a)$$

This cost corresponds to the following choice for the loss function:

$$Q(w; \mathbf{x}) = (d - \mathbf{u}w)^2 \quad (3b)$$

Such quadratic costs are widely used in estimation and adaptation problems [33]–[37]. They are also widely used as quadratic risk functions in machine learning applications [38], [39]. The gradient vector and Hessian matrix of $J(w)$ are easily seen to be:

$$\nabla_w J(w) = 2(R_u w - r_{du})^\top, \quad \nabla_w^2 J(w) = 2R_u \quad (3c)$$

◆

Example II.2 (Logistic or log-loss risks). Let γ denote a binary random variable that assumes the values ± 1 , and let \mathbf{h} denote an $M \times 1$ random (feature) vector with $R_h = \mathbb{E} \mathbf{h} \mathbf{h}^\top$. The combined quantities $\{\gamma, \mathbf{h}\}$ represent the random variable \mathbf{x} referred to in (2). In the context of machine learning and pattern classification problems [38]–[40], the variable γ designates the class that feature vector \mathbf{h} belongs to. In these problems, one seeks the vector w° that minimizes the regularized logistic risk function:

$$J(w) \triangleq \frac{\rho}{2} \|w\|^2 + \mathbb{E} \left\{ \ln \left[1 + e^{-\gamma \mathbf{h}^\top w} \right] \right\} \quad (4a)$$

where $\rho > 0$ is some regularization parameter, $\ln(\cdot)$ is the natural logarithm function, and $\|w\|^2 = w^\top w$. The risk (4a) corresponds to the following choice for the loss function:

$$Q(w; \mathbf{x}) \triangleq \frac{\rho}{2} \|w\|^2 + \ln \left[1 + e^{-\gamma \mathbf{h}^\top w} \right] \quad (4b)$$

Once w^o is recovered, its value can be used to classify new feature vectors, say, $\{\mathbf{h}_\ell\}$, into classes $+1$ or -1 . This can be achieved by assigning feature vectors with $\mathbf{h}_\ell^\top w^o \geq 0$ to one class and feature vectors with $\mathbf{h}_\ell^\top w^o < 0$ to another class. It can be easily verified that for the above $J(w)$:

$$\nabla_w J(w) = \rho w^\top - \mathbb{E} \left\{ \gamma \mathbf{h}^\top \cdot \frac{e^{-\gamma \mathbf{h}^\top w}}{1 + e^{-\gamma \mathbf{h}^\top w}} \right\} \quad (4c)$$

$$\nabla_w^2 J(w) = \rho I_M + \mathbb{E} \left\{ \mathbf{h} \mathbf{h}^\top \cdot \frac{e^{-\gamma \mathbf{h}^\top w}}{(1 + e^{-\gamma \mathbf{h}^\top w})^2} \right\} \quad (4d)$$

where I_M denotes the identity matrix of size $M \times M$.

B. Conditions on Cost Function

Stochastic gradient algorithms are powerful iterative procedures for solving optimization problems of the form

$$\underset{w}{\text{minimize}} \quad J(w) \quad (5)$$

While the analysis that follows can be pursued under more relaxed conditions (see, e.g., the treatments in [41]–[44]), it is sufficient for our purposes to require $J(w)$ to be strongly-convex and twice-differentiable with respect to w . The cost function $J(w)$ is said to be ν -strongly convex if, and only if, its Hessian matrix is sufficiently bounded away from zero [42], [45]–[47]:

$$J(w) \text{ is } \nu\text{-strongly convex} \iff \nabla_w^2 J(w) \geq \nu I_M > 0 \quad (6a)$$

for all w and for some scalar $\nu > 0$, where the notation $A > 0$ signifies that matrix A is positive-definite. Strong convexity is a useful condition in the context of adaptation and learning from streaming data because it helps guard against ill-conditioning in the algorithms; it also helps ensure that $J(w)$ has a *unique* global minimum, say, at location w^o ; there will be no other minima, maxima, or saddle points. In addition, it is well-known that strong convexity helps endow stochastic-gradient algorithms with geometric convergence rates in the order of $O(\alpha^i)$, for some $0 \leq \alpha < 1$ and where i is the iteration index [42], [43]. For comparison purposes, when the function $J(w)$ is convex but not necessarily strongly convex, then convexity is equivalent to the following condition:

$$J(w) \text{ is convex} \iff \nabla_w^2 J(w) \geq 0 \quad (6b)$$

for all w . In this case, there can now be multiple global minima. Moreover, the convergence of stochastic-gradient algorithms will occur at the slower rate of $O(1/i)$ [42], [43].

In many problems of interest in adaptation and learning, the cost function $J(w)$ is either already strongly convex or can be made strongly convex by means of regularization. For example, it is common in machine learning problems [38], [39] and in adaptation and estimation problems [35], [37] to incorporate regularization factors into the cost functions; these factors help ensure strong convexity. For instance, the mean-square-error cost (3a) is strongly convex whenever $R_u > 0$. If R_u happens to be singular, then the following regularized cost will be strongly convex:

$$J(w) \triangleq \frac{\rho}{2} \|w\|^2 + \mathbb{E}(\mathbf{d} - \mathbf{u}w)^2 \quad (7)$$

where $\rho > 0$ is a regularization parameter similar to (4a).

Besides strong convexity, we shall also assume that the gradient vector of $J(w)$ is δ -Lipschitz, namely, there exists $\delta > 0$ such that

$$\|\nabla_w J(w_2) - \nabla_w J(w_1)\| \leq \delta \|w_2 - w_1\| \quad (8)$$

for all w_1 and w_2 . It can be verified that for twice-differentiable costs, conditions (6a) and (8) combined are equivalent to

$$0 < \nu I_M \leq \nabla_w^2 J(w) \leq \delta I_M \quad (9)$$

For example, it is clear that the Hessian matrices in (3c) and (4d) satisfy this property since

$$2\lambda_{\min}(R_u)I_M \leq \nabla_w^2 J(w) \leq 2\lambda_{\max}(R_u)I_M \quad (10a)$$

in the first case and

$$\rho I_M \leq \nabla_w^2 J(w) \leq (\rho + \lambda_{\max}(R_h))I_M \quad (10b)$$

in the second case, where the notation $\lambda_{\min}(R)$ and $\lambda_{\max}(R)$ refers to the smallest and largest eigenvalues of the symmetric matrix argument, R , respectively. In summary, we will be assuming the following conditions [2], [3], [48], [49].

Assumption II.1 (Conditions on cost function). *The cost function $J(w)$ is twice-differentiable and satisfies (9) for some positive parameters $\nu \leq \delta$. Condition (9) is equivalent to requiring $J(w)$ to be ν -strongly convex and for its gradient vector to be δ -Lipschitz as in (6a) and (8), respectively. \square*

C. Stochastic-Gradient Approximation

The traditional gradient-descent algorithm for solving (5) takes the form:

$$w_i = w_{i-1} - \mu \nabla_{w^\top} J(w_{i-1}), \quad i \geq 0 \quad (11)$$

where $i \geq 0$ is an iteration index and $\mu > 0$ is a small step-size parameter. Starting from some initial condition, w_{-1} , the iterates $\{w_i\}$ correspond to successive estimates for the minimizer w^o . In order to run recursion (11), we need to have access to the true gradient vector. This information is generally unavailable in most instances involving learning from data. For example, when cost functions are defined as the expectations of certain loss functions as in (2), the statistical distribution of the data \mathbf{x} may not be known beforehand. In that case, the exact form of $J(w)$ will not be known since the expectation of $Q(w; \mathbf{x})$ cannot be computed. In such situations, it is customary to replace the true gradient vector, $\nabla_{w^\top} J(w_{i-1})$, by an instantaneous approximation for it, and which we shall denote by $\widehat{\nabla_{w^\top} J(w_{i-1})}$. Doing so leads to the following *stochastic-gradient* recursion in lieu of (11):

$$\mathbf{w}_i = \mathbf{w}_{i-1} - \mu \widehat{\nabla_{w^\top} J(w_{i-1})}, \quad i \geq 0 \quad (12)$$

We use the **boldface** notation, \mathbf{w}_i , for the iterates in (12) to highlight the fact that these iterates are now randomly perturbed versions of the values $\{w_i\}$ generated by the original recursion (11). The random perturbations arise from the use of the approximate gradient vector. The boldface notation

is therefore meant to emphasize the random nature of the iterates in (12). We refer to recursion (12) as a *synchronous* implementation since updates occur continuously over the iteration index i . This terminology is meant to distinguish the above recursion from its *asynchronous* counterpart, which is introduced and studied further ahead in Sec. II-E.

We illustrate construction (12) by considering a scenario from classical adaptive filter theory [33]–[35], where the gradient vector is approximated directly from data realizations. The construction will reveal why stochastic-gradient implementations of the form (12), using approximate rather than exact gradient information, become naturally endowed with the ability to respond to *streaming* data.

Example II.3 (LMS adaptation). Let $\mathbf{d}(i)$ denote a streaming sequence of zero-mean random variables with variance $\sigma_d^2 = \mathbb{E} \mathbf{d}^2(i)$. Let \mathbf{u}_i denote a streaming sequence of $1 \times M$ independent zero-mean random vectors with covariance matrix $R_u = \mathbb{E} \mathbf{u}_i^\top \mathbf{u}_i > 0$. Both processes $\{\mathbf{d}(i), \mathbf{u}_i\}$ are assumed to be jointly wide-sense stationary. The cross-covariance vector between $\mathbf{d}(i)$ and \mathbf{u}_i is denoted by $r_{du} = \mathbb{E} \mathbf{d}(i) \mathbf{u}_i^\top$. The data $\{\mathbf{d}(i), \mathbf{u}_i\}$ are assumed to be related via a linear regression model of the form:

$$\mathbf{d}(i) = \mathbf{u}_i w^\circ + v(i) \quad (13a)$$

for some unknown parameter vector w° , and where $v(i)$ is a zero-mean white-noise process with power $\sigma_v^2 = \mathbb{E} v^2(i)$ and assumed independent of \mathbf{u}_j for all i, j . Observe that we are using parentheses to represent the time-dependency of a scalar variable, such as writing $\mathbf{d}(i)$, and subscripts to represent the time-dependency of a vector variable, such as writing \mathbf{u}_i . This convention will be used throughout the chapter. In a manner similar to Example II.1, we again pose the problem of estimating w° by minimizing the mean-square error cost

$$J(w) = \mathbb{E} (\mathbf{d}(i) - \mathbf{u}_i w)^2 \equiv \mathbb{E} Q(w; \mathbf{x}_i) \quad (13b)$$

where now the quantities $\{\mathbf{d}(i), \mathbf{u}_i\}$ represent the random data \mathbf{x}_i in the definition of the loss function, $Q(w; \mathbf{x}_i)$. Using (11), the gradient-descent recursion in this case will take the form:

$$w_i = w_{i-1} - 2\mu [R_u w_{i-1} - r_{du}], \quad i \geq 0 \quad (13c)$$

The main difficulty in running this recursion is that it requires knowledge of the moments $\{r_{du}, R_u\}$. This information is rarely available beforehand; the adaptive agent senses instead realizations $\{\mathbf{d}(i), \mathbf{u}_i\}$ whose statistical distributions have moments $\{r_{du}, R_u\}$. The agent can use these realizations to approximate the moments and the true gradient vector. There are many constructions that can be used for this purpose, with different constructions leading to different adaptive algorithms [33]–[36]. It is sufficient to illustrate the construction by focusing on one of the most popular adaptive algorithms, which results from using the data $\{\mathbf{d}(i), \mathbf{u}_i\}$ to compute *instantaneous* approximations for the unavailable moments as follows:

$$r_{du} \approx \mathbf{d}(i) \mathbf{u}_i^\top, \quad R_u \approx \mathbf{u}_i^\top \mathbf{u}_i \quad (13d)$$

By doing so, the true gradient vector is approximated by:

$$\widehat{\nabla_{w^\top} J(w)} = 2 [\mathbf{u}_i^\top \mathbf{u}_i w - \mathbf{u}_i^\top \mathbf{d}(i)] = \nabla_{w^\top} Q(w; \mathbf{x}_i) \quad (13e)$$

Observe that this construction amounts to replacing the true gradient vector, $\nabla_{w^\top} J(w)$, by the gradient vector of the loss function itself (which, equivalently, amounts to dropping the expectation operator). Substituting (13e) into (13c) leads to the well-known (synchronous) least-mean-squares (LMS, for short) algorithm [33]–[35]:

$$w_i = w_{i-1} + 2\mu \mathbf{u}_i^\top [\mathbf{d}(i) - \mathbf{u}_i w_{i-1}], \quad i \geq 0 \quad (13f)$$

The LMS algorithm is therefore a stochastic-gradient algorithm. By relying directly on the instantaneous data $\{\mathbf{d}(i), \mathbf{u}_i\}$, the algorithm is infused with useful tracking abilities. This is because drifts in the model w° from (13a) will be reflected in the data $\{\mathbf{d}(i), \mathbf{u}_i\}$, which are used directly in the update (13f). \blacklozenge

The idea of using sample realizations to approximate actual expectations, as was the case with step (13e), is at the core of what is known as *stochastic approximation theory*. According to [35], [44], the pioneering work in the field of stochastic approximation is that of [50], which is a variation of a scheme developed about two decades earlier in [51]. The work by [50] dealt primarily with *scalar* weights w and was extended later by [52], [53] to weight *vectors* — see [54]. During the 1950s, stochastic approximation theory did not receive much attention in the engineering community until the landmark work by [55], which developed the real form of the LMS algorithm (13f). The algorithm has since then found remarkable success in a wide range of applications.

If desired, it is also possible to employ iteration-dependent step-size sequences, $\mu(i)$, in (12) instead of the constant step-size μ , and to require $\mu(i)$ to satisfy

$$\sum_{i=0}^{\infty} \mu^2(i) < \infty, \quad \sum_{i=0}^{\infty} \mu(i) = \infty \quad (14)$$

Under some technical conditions, it is well-known that such step-size sequences ensure the convergence of w_i towards w° almost surely as $i \rightarrow \infty$ [2], [42]–[44]. However, conditions (14) force the step-size sequence to decay to zero, which is problematic for applications requiring continuous adaptation and learning from streaming data. This is because, in such applications, it is not unusual for the location of the minimizer, w° , to drift with time. With $\mu(i)$ decaying towards zero, the stochastic-gradient algorithm (12) will stop updating and will not be able to track drifts in the solution. For this reason, we shall focus on constant step-sizes from this point onwards since we are interested in solutions with tracking abilities.

Now, the use of an approximate gradient vector in (12) introduces perturbations relative to the operation of the original recursion (11). We refer to the perturbation as gradient noise and define it as the difference:

$$s_i(w_{i-1}) \triangleq \widehat{\nabla_{w^\top} J(w_{i-1})} - \nabla_{w^\top} J(w_{i-1}) \quad (15)$$

The presence of this perturbation prevents the stochastic iterate, w_i , from converging almost-surely to the minimizer w° when constant step-sizes are used. Some deterioration in performance will occur and the iterate w_i will instead fluctuate close to w° . We will assess the size of these fluctuations by measuring their steady-state mean-square value (also called mean-square-deviation or MSD). It will turn out that the MSD is small and in the order of $O(\mu)$ — see (28c) further ahead. It will also turn out that stochastic-gradient algorithms converge towards their MSD levels at a geometric rate. In this way, we will be able to conclude that adaptation with small constant step-sizes can still lead to reliable performance in the presence of gradient noise,

which is a reassuring result. We will also be able to conclude that adaptation with constant step-sizes is useful even for stationary environments. This is because it is generally sufficient in practice to reach an iterate \mathbf{w}_i within some fidelity level from \mathbf{w}^o in a *finite* number of iterations. As long as the MSD level is satisfactory, a stochastic-gradient algorithm will be able to attain satisfactory fidelity within a reasonable time frame. In comparison, although diminishing step-sizes ensure almost-sure convergence of \mathbf{w}_i to \mathbf{w}^o , they nevertheless disable tracking and can only guarantee slower than geometric rates of convergence (see, e.g., [2], [2], [42], [43]). The next example from [48] illustrates the nature of the gradient noise process (15) in the context of mean-square-error adaptation.

Example II.4 (Gradient noise). It is clear from the expressions in Example II.3 that the corresponding gradient noise process is

$$\mathbf{s}_i(\mathbf{w}_{i-1}) = 2 \left(R_u - \mathbf{u}_i^\top \mathbf{u}_i \right) \tilde{\mathbf{w}}_{i-1} - 2\mathbf{u}_i^\top \mathbf{v}(i) \quad (16a)$$

where we introduced the error vector:

$$\tilde{\mathbf{w}}_i \triangleq \mathbf{w}^o - \mathbf{w}_i \quad (16b)$$

Let the symbol \mathcal{F}_{i-1} represent the collection of all possible random events generated by the past iterates $\{\mathbf{w}_j\}$ up to time $i-1$ (more formally, \mathcal{F}_{i-1} is the *filtration* generated by the random process \mathbf{w}_j for $j \leq i-1$):

$$\mathcal{F}_{i-1} \triangleq \text{filtration} \{ \mathbf{w}_{-1}, \mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_{i-1} \} \quad (16c)$$

It follows from the conditions on the random processes $\{\mathbf{u}_i, \mathbf{v}(i)\}$ in Example II.3 that

$$\mathbb{E} [\mathbf{s}_i(\mathbf{w}_{i-1}) | \mathcal{F}_{i-1}] = 0 \quad (16d)$$

$$\mathbb{E} [\|\mathbf{s}_i(\mathbf{w}_{i-1})\|^2 | \mathcal{F}_{i-1}] \leq 4c \|\tilde{\mathbf{w}}_{i-1}\|^2 + 4\sigma_v^2 \text{Tr}(R_u) \quad (16e)$$

for some constant $c \geq 0$. If we take expectations of both sides of (16e), we further conclude that the variance of the gradient noise, $\mathbb{E} \|\mathbf{s}_i(\mathbf{w}_{i-1})\|^2$, is bounded by the combination of two factors. The first factor depends on the quality of the iterate, $\mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2$, while the second factor depends on σ_v^2 . Therefore, even if the adaptive agent is able to approach \mathbf{w}^o with great fidelity so that $\mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2$ is small, the size of the gradient noise will still depend on σ_v^2 . ♦

D. Conditions on Gradient Noise Process

In order to examine the convergence and performance properties of the stochastic-gradient recursion (12), it is necessary to introduce some assumptions on the stochastic nature of the gradient noise process, $\mathbf{s}_i(\cdot)$. The conditions that we introduce in the sequel are similar to conditions used earlier in the optimization literature, e.g., in [42, pp. 95–102] and [56, p. 635]; they are also motivated by the conditions we observed in the mean-square-error case in Example II.4. Following the developments in [3], [48], [49], we let

$$R_{s,i}(\mathbf{w}_{i-1}) \triangleq \mathbb{E} [\mathbf{s}_i(\mathbf{w}_{i-1}) \mathbf{s}_i^\top(\mathbf{w}_{i-1}) | \mathcal{F}_{i-1}] \quad (17a)$$

denote the conditional second-order moment of the gradient noise process, which generally depends on i . We assume that, in the limit, the covariance matrix tends to a constant value

when evaluated at \mathbf{w}^o and is denoted by

$$R_s \triangleq \lim_{i \rightarrow \infty} \mathbb{E} [\mathbf{s}_i(\mathbf{w}^o) \mathbf{s}_i^\top(\mathbf{w}^o) | \mathcal{F}_{i-1}] \quad (17b)$$

For example, comparing with expression (16a) for mean-square-error costs, we have

$$\mathbf{s}_i(\mathbf{w}^o) = -2\mathbf{u}_i^\top \mathbf{v}(i) \quad (18a)$$

$$R_s = 4\sigma_v^2 R_u \quad (18b)$$

Assumption II.2 (Conditions on gradient noise). It is assumed that the first and second-order conditional moments of the gradient noise process satisfy (17b) and

$$\mathbb{E} [\mathbf{s}_i(\mathbf{w}_{i-1}) | \mathcal{F}_{i-1}] = 0 \quad (19a)$$

$$\mathbb{E} [\|\mathbf{s}_i(\mathbf{w}_{i-1})\|^2 | \mathcal{F}_{i-1}] \leq \beta^2 \|\tilde{\mathbf{w}}_{i-1}\|^2 + \sigma_s^2 \quad (19b)$$

almost surely, for some nonnegative scalars β^2 and σ_s^2 . □

Condition (19a) ensures that the approximate gradient vector is unbiased. It follows from conditions (19a)–(19b) that the gradient noise process itself satisfies:

$$\mathbb{E} \mathbf{s}_i(\mathbf{w}_{i-1}) = 0 \quad (20a)$$

$$\mathbb{E} \|\mathbf{s}_i(\mathbf{w}_{i-1})\|^2 \leq \beta^2 \mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2 + \sigma_s^2 \quad (20b)$$

It is straightforward to verify that the gradient noise process (16a) in the mean-square-error case satisfies conditions (19a)–(19b). Note in particular from (16e) that we can make the identifications $\sigma_s^2 \rightarrow 4\sigma_v^2 \text{Tr}(R_u)$ and $\beta^2 \rightarrow 4c$.

E. Random Updates

We examined the performance of synchronous updates of the form (12) in some detail in [1], [2]. As indicated earlier, the focus of the current chapter is on extending the treatment from [1] to asynchronous implementations. Accordingly, the first main digression in the exposition relative to [1] occurs at this stage.

Thus, note that the stochastic-gradient recursion (12) employs a constant step-size parameter, $\mu > 0$. This means that this implementation expects the approximate gradient vector, $\widehat{\nabla_{\mathbf{w}^\top} J(\mathbf{w}_{i-1})}$, to be available at every iteration. Nevertheless, there are situations where data may arrive at the agent at random times, in which case the updates will also be occurring at random times. One way to capture this behavior is to model the step-size parameter as a *random* process, which we shall denote by the boldface notation $\boldsymbol{\mu}(i)$ (since boldface letters in our notation refer to random quantities). Doing so, allows us to replace the synchronous implementation (12) by the following asynchronous recursion:

$$\mathbf{w}_i = \mathbf{w}_{i-1} - \boldsymbol{\mu}(i) \widehat{\nabla_{\mathbf{w}^\top} J(\mathbf{w}_{i-1})}, \quad i \geq 0 \quad (21)$$

Observe that we are attaching a time index to the step-size parameter in order to highlight that its value will now be changing randomly from one iteration to another.

For example, one particular instance discussed further ahead in Example II.5 is when $\mu(i)$ is a Bernoulli random variable assuming one of two possible values, say, $\mu(i) = \mu > 0$ with probability p_μ and $\mu(i) = 0$ with probability $1 - p_\mu$. In this case, recursion (21) will be updating p_μ —fraction of the time. We will not limit our presentation to the Bernoulli model but will allow the random step-size to assume broader probability distributions; we denote its first and second-order moments as follows.

Assumption II.3 (Conditions on step-size process). *It is assumed that the stochastic process $\{\mu(i), i \geq 0\}$ consists of a sequence of independent and bounded random variables, $\mu(i) \in [0, \mu_{\text{ub}}]$, where $\mu_{\text{ub}} > 0$ is a constant upper bound. The mean and variance of $\mu(i)$ are fixed over i , and they are denoted by:*

$$\bar{\mu} \triangleq \mathbb{E} \mu(i) \quad (22a)$$

$$\sigma_\mu^2 \triangleq \mathbb{E} (\mu(i) - \bar{\mu})^2 \quad (22b)$$

with $\bar{\mu} > 0$ and $\sigma_\mu^2 \geq 0$. Moreover, it is assumed that, for any i , the random variable $\mu(i)$ is independent of any other random variable in the learning algorithm. \square

The following variable will play an important role in characterizing the mean-square-error performance of asynchronous updates and, hence, we introduce a unique symbol for it:

$$\mu_x \triangleq \bar{\mu} + \frac{\sigma_\mu^2}{\bar{\mu}} \quad (23)$$

where the subscript “ x ” is meant to refer to the “asynchronous” mode of operation. This expression captures the first and second-order moments of the random variations in the step-size parameter into a single variable, μ_x . While the constant step-size μ determines the performance of the synchronous implementation (12), it turns out that the constant variable μ_x defined above will play an equivalent role for the asynchronous implementation (21). Note further that the synchronous stochastic-gradient iteration (12) can be viewed as a special case of recursion (21) when the variance of $\mu(i)$ is set to zero, i.e., $\sigma_\mu^2 = 0$, and the mean value $\bar{\mu}$ is set to μ . Therefore, by using these substitutions, we will be able to deduce performance metrics for (12) from the performance metrics that we shall present for (21).

The following two examples illustrate situations involving random updates.

Example II.5 (Random updates under a Bernoulli model). Assume that at every iteration i , the agent adopts a random “on-off” policy to reduce energy consumption. It tosses a coin to decide whether to enter an active learning mode or a sleeping mode. Let $0 < p_\mu < 1$ denote the probability of entering the active mode. During the active mode, the agent employs a step-size value μ . This model is useful for situations in which data arrives randomly at the agent: at every iteration i , new data is available with probability p_μ . The random step-size process is therefore of the following type:

$$\mu(i) = \begin{cases} \mu, & \text{with probability } p_\mu \\ 0, & \text{with probability } 1 - p_\mu \end{cases} \quad (24a)$$

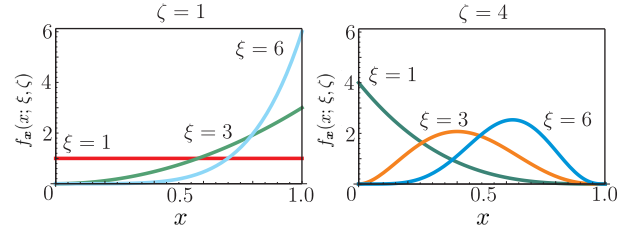


Fig. 1. The pdf of the Beta distribution, $f_x(x; \xi, \zeta)$, defined by (25a) for different values of the shape parameters ξ and ζ .

In this case, the mean and variance of $\mu(i)$ are given by:

$$\bar{\mu} = p_\mu \mu, \quad \sigma_\mu^2 = p_\mu(1 - p_\mu)\mu^2, \quad \mu_x = \mu \quad (24b)$$

◆

Example II.6 (Random updates under a Beta model). Since the random step-size $\mu(i)$ is limited to a finite-length interval $[0, \mu_{\text{ub}}]$, we may extend the Bernoulli model from the previous example by adopting a more general continuous Beta distribution for $\mu(i)$. The Beta distribution is an extension of the Bernoulli distribution. While the Bernoulli distribution assumes two discrete possibilities for the random variable, say, $\{0, \mu\}$, the Beta distribution allows for any value in the continuum $[0, \mu]$.

Thus, let x denote a generic scalar random variable that assumes values in the interval $[0, 1]$ according to a Beta distribution. Then, according to this distribution, the pdf of x , denoted by $f_x(x; \xi, \zeta)$, is determined by two shape parameters $\{\xi, \zeta\}$ as follows [57], [58]:

$$f_x(x; \xi, \zeta) = \begin{cases} \frac{\Gamma(\xi + \zeta)}{\Gamma(\xi)\Gamma(\zeta)} x^{\xi-1} (1-x)^{\zeta-1}, & 0 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases} \quad (25a)$$

where $\Gamma(\cdot)$ denotes the Gamma function [59], [60]. Figure 1 plots $f_x(x; \xi, \zeta)$ for two values of ζ . The mean and variance of the Beta distribution (25a) are given by:

$$\bar{x} = \frac{\xi}{\xi + \zeta}, \quad \sigma_x^2 = \frac{\xi\zeta}{(\xi + \zeta)^2(\xi + \zeta + 1)} \quad (25b)$$

We note that the classical uniform distribution over the interval $[0, 1]$ is a special case of the Beta distribution for $\xi = \zeta = 1$ — see Fig. 1. Likewise, the Bernoulli distribution with $p_\mu = 1/2$ is recovered from the Beta distribution by letting $\xi = \zeta \rightarrow 0$.

In the Beta model for asynchronous adaptation, we assume that the ratio $\mu(i)/\mu_{\text{ub}}$ follows a Beta distribution with parameters $\{\xi, \zeta\}$. Under this model, the mean and variance of the random step-size become:

$$\bar{\mu} = \left(\frac{\xi}{\xi + \zeta} \right) \mu_{\text{ub}}, \quad \sigma_\mu^2 = \left(\frac{\xi\zeta}{(\xi + \zeta)^2(\xi + \zeta + 1)} \right) \mu_{\text{ub}}^2 \quad (26)$$

◆

F. Mean-Square-Error Stability

We now examine the convergence of the asynchronous stochastic-gradient recursion (21). In the statement below, the notation $a = O(\mu)$ means $a \leq b\mu$ for some constant b that is independent of μ .

Lemma II.1 (Mean-square-error stability). *Assume the conditions under Assumptions II.1, II.2, and II.3 on the cost function, the gradient noise process, and the random step-size process hold. Let $\mu_o = 2\nu/(\delta^2 + \beta^2)$. For any μ_x satisfying*

$$\mu_x < \mu_o \quad (27)$$

it holds that $\mathbb{E} \|\tilde{\mathbf{w}}_i\|^2$ converges exponentially (i.e., at a geometric rate) according to the recursion

$$\mathbb{E} \|\tilde{\mathbf{w}}_i\|^2 \leq \alpha \mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2 + (\bar{\mu}^2 + \sigma_\mu^2) \sigma_s^2 \quad (28a)$$

where the scalar α satisfies $0 \leq \alpha < 1$ and is given by

$$\begin{aligned} \alpha &\triangleq 1 - 2\nu\bar{\mu} + (\delta^2 + \beta^2)(\bar{\mu}^2 + \sigma_\mu^2) \\ &= 1 - 2\nu\bar{\mu} + O(\mu_x^2) \end{aligned} \quad (28b)$$

It follows from (28a) that, for sufficiently small step-sizes:

$$\limsup_{i \rightarrow \infty} \mathbb{E} \|\tilde{\mathbf{w}}_i\|^2 = O(\mu_x) \quad (28c)$$

Proof: We subtract w^o from both sides of (21) to get

$$\tilde{\mathbf{w}}_i = \tilde{\mathbf{w}}_{i-1} + \boldsymbol{\mu}(i) \nabla_{w^\top} J(\mathbf{w}_{i-1}) + \boldsymbol{\mu}(i) \mathbf{s}_i(\mathbf{w}_{i-1}) \quad (29a)$$

We now appeal to the mean-value theorem [2], [42], [61] to write:

$$\begin{aligned} \nabla_{w^\top} J(\mathbf{w}_{i-1}) &= - \left[\int_0^1 \nabla_w^2 J(w_o - t\tilde{\mathbf{w}}_{i-1}) dt \right] \tilde{\mathbf{w}}_{i-1} \\ &\triangleq -\mathbf{H}_{i-1} \tilde{\mathbf{w}}_{i-1} \end{aligned} \quad (29b)$$

where we are introducing the *symmetric* and *random* time-variant matrix \mathbf{H}_{i-1} to represent the integral expression. Substituting into (29a), we get

$$\tilde{\mathbf{w}}_i = [I_M - \boldsymbol{\mu}(i) \mathbf{H}_{i-1}] \tilde{\mathbf{w}}_{i-1} + \boldsymbol{\mu}(i) \mathbf{s}_i(\mathbf{w}_{i-1}) \quad (29c)$$

so that from Assumption II.2:

$$\begin{aligned} \mathbb{E} [\|\tilde{\mathbf{w}}_i\|^2 | \mathcal{F}_{i-1}] &\leq (\mathbb{E} [\|I_M - \boldsymbol{\mu}(i) \mathbf{H}_{i-1}\|^2 | \mathcal{F}_{i-1}]) \|\tilde{\mathbf{w}}_{i-1}\|^2 + \\ &\quad (\mathbb{E} \boldsymbol{\mu}^2(i)) (\mathbb{E} [\|\mathbf{s}_i(\mathbf{w}_{i-1})\|^2 | \mathcal{F}_{i-1}]) \end{aligned} \quad (30)$$

It follows from (9) that

$$\begin{aligned} \|I_M - \boldsymbol{\mu}(i) \mathbf{H}_{i-1}\|^2 &= [\rho(I_M - \boldsymbol{\mu}(i) \mathbf{H}_{i-1})]^2 \\ &\leq \max \{ [1 - \boldsymbol{\mu}(i) \delta]^2, [1 - \boldsymbol{\mu}(i) \nu]^2 \} \\ &\leq 1 - 2\boldsymbol{\mu}(i) \nu + \boldsymbol{\mu}^2(i) \delta^2 \end{aligned} \quad (31a)$$

since $\nu \leq \delta$. In the first line above, the notation $\rho(A)$ denotes the spectral radius of its matrix argument (i.e., $\rho(A) = \max_k |\lambda_k(A)|$ in terms of the largest magnitude eigenvalue of A). From (31a) we obtain

$$\mathbb{E} (\|I_M - \boldsymbol{\mu}(i) \mathbf{H}_{i-1}\|^2 | \mathcal{F}_{i-1}) \leq 1 - 2\bar{\mu}\nu + (\bar{\mu}^2 + \sigma_\mu^2) \delta^2 \quad (31b)$$

Taking expectations of both sides of (30), we arrive at (28a) from (20b) and (31b) with α given by (28b). The bound in (27) on the moments of the random step-size ensures that $0 \leq \alpha < 1$. For the $O(\mu_x^2)$ approximation in expression (28b) note from (23) that

$$(\bar{\mu}^2 + \sigma_\mu^2) = \bar{\mu}\mu_x \leq \mu_x^2 \quad (32)$$

Iterating recursion (28a) gives

$$\mathbb{E} \|\tilde{\mathbf{w}}_i\|^2 \leq \alpha^{i+1} \mathbb{E} \|\tilde{\mathbf{w}}_{-1}\|^2 + \frac{(\bar{\mu}^2 + \sigma_\mu^2) \sigma_s^2}{1 - \alpha} \quad (33a)$$

Since $0 \leq \alpha < 1$, there exists an iteration value I_o large enough such that

$$\alpha^{i+1} \mathbb{E} \|\tilde{\mathbf{w}}_{-1}\|^2 \leq \frac{(\bar{\mu}^2 + \sigma_\mu^2) \sigma_s^2}{1 - \alpha}, \quad i > I_o \quad (33b)$$

It follows that the variance $\mathbb{E} \|\tilde{\mathbf{w}}_i\|^2$ converges exponentially to a region that is upper bounded by $2(\bar{\mu}^2 + \sigma_\mu^2) \sigma_s^2 / (1 - \alpha)$. It can be verified that this bound does not exceed $2\mu_x \sigma_s^2 / \nu$, which is $O(\mu_x)$, for any $\mu_x < \mu_o/2$.

G. Mean-Square-Error Performance

We conclude from (28c) that the mean-square error can be made as small as desired by using small step-sizes, μ_x . In this section we derive a closed-form expression for the asymptotic mean-square error, which is more frequently called the *mean-square deviation* (MSD) and is defined as:

$$\text{MSD} \triangleq \lim_{i \rightarrow \infty} \mathbb{E} \|\tilde{\mathbf{w}}_i\|^2 \quad (34a)$$

Strictly speaking, the limit on the right-hand side of the above expression may not exist. A more accurate definition for the MSD appears in Eq. (4.86) of [2], namely,

$$\text{MSD} \triangleq \mu_x \cdot \left(\lim_{\mu_x \rightarrow 0} \limsup_{i \rightarrow \infty} \frac{1}{\mu_x} \mathbb{E} \|\tilde{\mathbf{w}}_i\|^2 \right) \quad (34b)$$

However, it was explained in [2][Sec. 4.5] that derivations that assume the validity of (34a) still lead to the same expression for the MSD to first-order in μ_x as derivations that rely on the more formal definition (34b). We therefore continue with (34a) for simplicity of presentation in this chapter.

Evaluating the MSD of a learning algorithm is a useful step to pursue because once performance expressions are available, it becomes possible to compare different configurations for adaptation and learning (such as non-cooperative, centralized, and distributed implementations). It also becomes possible to quantify how performance depends on the algorithm and system parameters (such as step-size, network topology, and cooperation policy); these parameters can then be optimized for enhanced performance. We explain below how an expression for the MSD can be obtained by following the energy conservation technique of [2], [35], [36], [62], [63]. For that purpose, we need to introduce two smoothness conditions.

Assumption II.4 (Smoothness conditions). *In addition to Assumptions II.1 and II.2, we assume that the Hessian matrix of the cost function and the noise covariance matrix defined by (17a) are locally Lipschitz continuous in a small neighborhood around $w = w^o$:*

$$\|\nabla_w^2 J(w^o + \delta w) - \nabla_w^2 J(w^o)\| \leq \tau \|\delta w\| \quad (35a)$$

$$\|R_{s,i}(w^o + \delta w) - R_{s,i}(w^o)\| \leq \tau_2 \|\delta w\|^\kappa \quad (35b)$$

for small perturbations $\|\delta w\| \leq r$ and for some $\tau, \tau_2 \geq 0$ and $1 \leq \kappa \leq 2$. \square

The range of values for κ can be enlarged, e.g., to $\kappa \in (0, 4]$. The only change in allowing a wider range for κ is that the exponent of the higher-order term, $O(\mu_x^{3/2})$, that will appear in several performance expressions, as is the case with (41a)–(41b), will need to be adjusted from $\frac{3}{2}$ to $\min\{\frac{3}{2}, 1 + \frac{\kappa}{2}\}$, without affecting the first-order term that determines the MSD [2], [4], [64]. Therefore, it is sufficient to continue with $\kappa \in [1, 2]$ to illustrate the key concepts though the MSD expressions will still be valid to first-order in μ_x .

Using (9), it can be verified that condition (35a) translates into a global Lipschitz property relative to the minimizer w^o , i.e., it will also hold that [2], [4]:

$$\|\nabla_w^2 J(w) - \nabla_w^2 J(w^o)\| \leq \tau' \|w - w^o\| \quad (35c)$$

for all w and for some $\tau' \geq 0$. For example, both conditions (35a)–(35b) are readily satisfied by mean-square-error costs. Using property (35c), we can now motivate a useful long-term model for the evolution of the error vector \tilde{w}_i after sufficient iterations, i.e., for $i \gg 1$. Indeed, let us reconsider recursion (29c) and introduce the deviation matrix:

$$\tilde{H}_{i-1} \triangleq H - H_{i-1} \quad (36a)$$

where the constant (symmetric and positive-definite) matrix H is defined as:

$$H \triangleq \nabla_w^2 J(w^o) \quad (36b)$$

Substituting (36a) into (29c) gives

$$\tilde{w}_i = (I_M - \mu(i) H) \tilde{w}_{i-1} + \mu(i) s_i(w_{i-1}) + \mu(i) c_{i-1} \quad (37a)$$

where

$$c_{i-1} \triangleq \tilde{H}_{i-1} \tilde{w}_{i-1} \quad (37b)$$

Using (35c) and the fact that $(\mathbb{E} a)^2 \leq \mathbb{E} a^2$ for any real-valued random variable, we can bound the conditional expectation of the norm of the perturbation term as follows:

$$\begin{aligned} \mathbb{E} [\|\mu(i) c_{i-1}\| | \mathcal{F}_{i-1}] &= (\mathbb{E} \mu(i)) (\mathbb{E} [\|c_{i-1}\| | \mathcal{F}_{i-1}]) \\ &\leq \sqrt{(\mathbb{E} \mu^2(i))} \mathbb{E} [\|c_{i-1}\| | \mathcal{F}_{i-1}] \\ &\stackrel{(35c)}{=} \sqrt{\bar{\mu}^2 + \sigma_\mu^2} \cdot \frac{\tau'}{2} \|\tilde{w}_{i-1}\|^2 \\ &\leq \frac{\bar{\mu}^2 + \sigma_\mu^2}{\bar{\mu}} \cdot \frac{\tau'}{2} \|\tilde{w}_{i-1}\|^2 \\ &= \mu_x \cdot \frac{\tau'}{2} \|\tilde{w}_{i-1}\|^2 \end{aligned} \quad (38a)$$

so that using (28c), we conclude that:

$$\limsup_{i \rightarrow \infty} \mathbb{E} \|\mu(i) c_{i-1}\| = O(\mu_x^2) \quad (38b)$$

We can deduce from this result that $\|\mu(i) c_{i-1}\| = O(\mu_x^2)$ asymptotically with *high probability* [2], [4]. To see this, let $r_c = m\mu_x^2$, for any constant integer $m \geq 1$. Now, calling upon Markov's inequality [65]–[67], we conclude from (38b) that for $i \gg 1$:

$$\begin{aligned} \Pr(\|\mu(i) c_{i-1}\| < r_c) &= 1 - \Pr(\|\mu(i) c_{i-1}\| \geq r_c) \\ &\geq 1 - \frac{\mathbb{E} \|\mu(i) c_{i-1}\|}{r_c} \\ &\stackrel{(38b)}{\geq} 1 - O(1/m) \end{aligned} \quad (38c)$$

This result shows that the probability of having $\|\mu(i) c_{i-1}\|$ bounded by r_c can be made arbitrarily close to one by selecting a large enough value for m . Once the value for m has been fixed to meet a desired confidence level, then $r_c = O(\mu_x^2)$. This analysis, along with recursion (37a), motivate us to assess the mean-square performance of the error recursion (29c) by considering instead the following long-term model, which holds with high probability after sufficient iterations $i \gg 1$:

$$\tilde{w}_i = (I_M - \mu(i) H) \tilde{w}_{i-1} + \mu(i) s_i(w_{i-1}) + O(\mu_x^2) \quad (39)$$

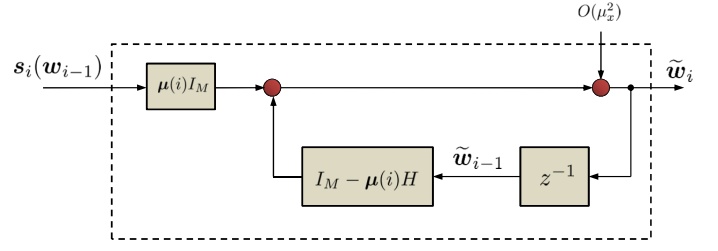


Fig. 2. A block-diagram representation of the long-term recursion (39) for single-agent adaptation and learning.

Working with iteration (39) is helpful because its dynamics is driven by the constant matrix H as opposed to the random matrix H_{i-1} in the original error recursion (29c). If desired, it can be shown that, under some technical conditions on the fourth-order moment of the gradient noise process, the MSD expression that will result from using (39) is within $O(\mu_x^{3/2})$ of the actual MSD expression for the original recursion (29c) — see [2], [4], [64] for a formal proof of this fact. Therefore, it is sufficient to rely on the long-term model (39) to obtain performance expressions that are accurate to first-order in μ_x . Figure 2 provides a block-diagram representation for (39).

Before explaining how model (39) can be used to assess the MSD, we remark that there is a second useful metric for evaluating the performance of stochastic gradient algorithms. This metric relates to the mean excess-cost; which is also called the *excess-risk* (ER) in the machine learning literature [38], [39] and the excess-mean-square-error (EMSE) in the adaptive filtering literature [33]–[35]. We denote it by the letters ER and define it as the average fluctuation of the cost function around its minimum value:

$$\text{ER} \triangleq \lim_{i \rightarrow \infty} \mathbb{E} \{J(w_{i-1}) - J(w^o)\} \quad (40a)$$

Using the smoothness condition (35a), and the mean-value theorem [42], [61] again, it can be verified that [2], [4], [64]:

$$\text{ER} \triangleq \lim_{i \rightarrow \infty} \mathbb{E} \|\tilde{w}_{i-1}\|_{\frac{1}{2}H}^2 + O(\mu_x^{3/2}) \quad (40b)$$

Lemma II.2 (Mean-square-error performance). *Assume the conditions under Assumptions II.1, II.2, II.3, and II.4 on the cost function, the gradient noise process, and the random step-size process hold. Assume further that the asynchronous step-size parameter μ_x is sufficiently small to ensure mean-square stability as required by (27). Then, the MSD and ER metrics for the asynchronous stochastic-gradient algorithm (21) are well-approximated to first-order in μ_x by the expressions:*

$$\text{MSD}^{\text{asyn}} = \frac{\mu_x}{2} \text{Tr}(H^{-1} R_s) + O(\mu_x^{3/2}) \quad (41a)$$

$$\text{ER}^{\text{asyn}} = \frac{\mu_x}{4} \text{Tr}(R_s) + O(\mu_x^{3/2}) \quad (41b)$$

where R_s and H are defined by (17b) and (36b), and where we are adding the superscript “asyn” for clarity in order to distinguish these measures from the corresponding measures in the synchronous case (mentioned below in (43a)–(43c)). Moreover, we derived earlier in (28b) the following expression for the convergence rate:

$$\alpha^{\text{asyn}} = 1 - 2\nu\bar{\mu} + O(\mu_x^2) \quad (41c)$$

Proof: We introduce the eigen-decomposition $H = U\Lambda U^\top$, where U is orthonormal and Λ is diagonal with positive entries, and rewrite (39) in terms of transformed quantities:

$$\bar{\mathbf{w}}_i = (I_M - \mu(i)\Lambda)\bar{\mathbf{w}}_{i-1} + \mu(i)\bar{\mathbf{s}}_i(\mathbf{w}_{i-1}) + O(\mu_x^2) \quad (42a)$$

where $\bar{\mathbf{w}}_i = U^\top \tilde{\mathbf{w}}_i$ and $\bar{\mathbf{s}}_i(\mathbf{w}_{i-1}) = U^\top \mathbf{s}_i(\mathbf{w}_{i-1})$. Let Σ denote an arbitrary $M \times M$ diagonal matrix with positive entries that we are free to choose. Then, equating the weighted squared norms of both sides of (42a) and taking expectations gives for $i \gg 1$:

$$\mathbb{E} \|\bar{\mathbf{w}}_i\|_\Sigma^2 = \mathbb{E} \|\bar{\mathbf{w}}_{i-1}\|_\Sigma^2 + (\bar{\mu}^2 + \sigma_\mu^2) \mathbb{E} \|\bar{\mathbf{s}}_i(\mathbf{w}_{i-1})\|_\Sigma^2 + O(\mu_x^{5/2}) \quad (42b)$$

where

$$\Sigma' \triangleq \mathbb{E} (I_M - \mu(i)\Lambda)\Sigma(I_M - \mu(i)\Lambda) = \Sigma - 2\bar{\mu}\Lambda\Sigma + O(\mu_x^2) \quad (42c)$$

From (17b), (19b), (28c), and (35b) we obtain:

$$\lim_{i \rightarrow \infty} \mathbb{E} \|\bar{\mathbf{s}}_i(\mathbf{w}_{i-1})\|_\Sigma^2 = \text{Tr}(U\Sigma U^\top R_s) + O(\mu_x^{\kappa/2}) \quad (42d)$$

Therefore, substituting into (42b) gives for $i \rightarrow \infty$:

$$\lim_{i \rightarrow \infty} \mathbb{E} \|\bar{\mathbf{w}}_i\|_{2\Lambda\Sigma}^2 = \mu_x \text{Tr}(U\Sigma U^\top R_s) + O(\mu_x^{3/2}) \quad (42e)$$

Since we are free to choose Σ , we let $\Sigma = \frac{1}{2}\Lambda^{-1}$ and arrive at (41a) since $\|\bar{\mathbf{w}}_i\|^2 = \|\tilde{\mathbf{w}}_i\|^2$ and $U\Sigma U^\top = \frac{1}{2}H^{-1}$. On the other hand, selecting $\Sigma = \frac{1}{4}I_M$ leads to (41b). \blacksquare

We recall our earlier remark that the synchronous stochastic-gradient recursion (12) can be viewed as a special case of the asynchronous update (21) by setting $\sigma_\mu^2 = 0$ and $\mu_x = \bar{\mu} \equiv \mu$. Substituting these values into (41a)–(41c), we obtain for the synchronous implementation (12):

$$\text{MSD}^{\text{sync}} = \frac{\mu}{2} \text{Tr}(H^{-1}R_s) + O(\mu^{3/2}) \quad (43a)$$

$$\text{ER}^{\text{sync}} = \frac{\mu}{2} \text{Tr}(R_s) + O(\mu^{3/2}) \quad (43b)$$

$$\alpha^{\text{sync}} = 1 - 2\nu\mu + O(\mu^2) \quad (43c)$$

which are the same expressions presented in [1] and which agree with classical results for LMS adaptation [68]–[73]. The matrices R_s that appear in these expressions, and in (41a)–(41c), were defined earlier in (17b) and they correspond to the covariance matrices of the gradient noise processes in their respective (synchronous or asynchronous) implementations.

The examples that follow show how expressions (41a) and (41b) can be used to recover performance metrics for mean-square-error adaptation and learning under random updates.

Example II.7 (Performance of asynchronous LMS adaptation). We reconsider the LMS recursion (13f) albeit in an asynchronous mode of operation, namely,

$$\mathbf{w}_i = \mathbf{w}_{i-1} + 2\mu(i)\mathbf{u}_i^\top [\mathbf{d}(i) - \mathbf{u}_i\mathbf{w}_{i-1}], \quad i \geq 0 \quad (44)$$

We know from Example II.3 and (18a) that this situation corresponds to $H = 2R_u$ and $R_s = 4\sigma_v^2 R_u$. Substituting into (41a) and (41b) leads to the following expressions for the MSD and EMSE of the asynchronous LMS filter:

$$\text{MSD}_{\text{LMS}}^{\text{asyn}} \approx \mu_x M \sigma_v^2 = O(\mu_x) \quad (45a)$$

$$\text{EMSE}_{\text{LMS}}^{\text{asyn}} \approx \mu_x \sigma_v^2 \text{Tr}(R_u) = O(\mu_x) \quad (45b)$$

where here, and elsewhere, we will be using the notation \approx to indicate that we are ignoring higher-order terms in μ_x .

For example, let us assume a Bernoulli update model for $\mu(i)$ where the filter updates with probability p_μ using a step-size value μ or stays inactive otherwise. In this case, we conclude from (24b) that $\mu_x = \mu$ so that the above performance expressions for the MSD and EMSE metrics will coincide with the values obtained in the synchronous case as well. In other words, as long as $p_\mu > 0$, and given sufficient time to learn, the steady-state performance levels are not affected whether the algorithm learns in a synchronous or asynchronous manner. However, the convergence rate is affected since $\bar{\mu} = \mu p_\mu$ and, therefore,

$$\alpha_{\text{LMS}}^{\text{sync}} \approx 1 - 2\nu\mu \quad (46a)$$

$$\alpha_{\text{LMS}}^{\text{asyn}} \approx 1 - 2\nu\mu p_\mu > \alpha_{\text{LMS}}^{\text{sync}} \quad (46b)$$

It follows that asynchronous LMS adaptation attains the same performance levels as synchronous LMS adaptation albeit at a slower convergence rate.

We may alternatively compare the performance of the synchronous and asynchronous implementations by fixing their convergence rates to the same value. Thus, consider now a second random update scheme with mean $\bar{\mu}$ and let us set $\mu = \bar{\mu}$. That is, the step-size used by the synchronous implementation is set equal to the mean step-size used by the asynchronous implementation. Then, in this case, we will get $\alpha_{\text{LMS}}^{\text{sync}} = \alpha_{\text{LMS}}^{\text{asyn}}$ so that the convergence rates coincide to first-order. However, it now holds that $\text{MSD}_{\text{LMS}}^{\text{asyn}} > \text{MSD}_{\text{LMS}}^{\text{sync}}$ since $\bar{\mu}_x > \mu$ so that some deterioration in MSD performance occurs in the asynchronous environment. \blacklozenge

Example II.8 (Performance of asynchronous online learners). Consider a stand-alone learner receiving a streaming sequence of independent data vectors $\{\mathbf{x}_i, i \geq 0\}$ that arise from some fixed probability distribution \mathcal{X} . The goal is to learn the vector \mathbf{w}^o that optimizes some ν -strongly convex risk function $J(\mathbf{w})$ defined in terms of a loss function [75], [76]:

$$\mathbf{w}^o \triangleq \arg \min_{\mathbf{w}} J(\mathbf{w}) = \arg \min_{\mathbf{w}} \mathbb{E} Q(\mathbf{w}; \mathbf{x}_i) \quad (47a)$$

In an asynchronous environment, the learner seeks \mathbf{w}^o by running the stochastic-gradient algorithm with random step-sizes:

$$\mathbf{w}_i = \mathbf{w}_{i-1} - \mu(i) \nabla_{\mathbf{w}^\top} Q(\mathbf{w}_{i-1}; \mathbf{x}_i), \quad i \geq 0 \quad (47b)$$

The gradient noise vector is still given by

$$\mathbf{s}_i(\mathbf{w}_{i-1}) = \nabla_{\mathbf{w}^\top} Q(\mathbf{w}_{i-1}; \mathbf{x}_i) - \nabla_{\mathbf{w}^\top} J(\mathbf{w}_{i-1}) \quad (47c)$$

Since $\nabla_{\mathbf{w}} J(\mathbf{w}^o) = 0$, and since the distribution of \mathbf{x}_i is stationary, it follows that the covariance matrix of $\mathbf{s}_i(\mathbf{w}^o)$ is constant and given by

$$R_s = \mathbb{E} \nabla_{\mathbf{w}^\top} Q(\mathbf{w}^o; \mathbf{x}_i) \nabla_{\mathbf{w}} Q(\mathbf{w}^o; \mathbf{x}_i) \quad (47d)$$

The excess-risk measure that will result from this stochastic implementation is then given by (41b) so that $\text{ER} = O(\mu_x)$. \blacklozenge

III. CENTRALIZED ADAPTATION AND LEARNING

The discussion in the previous section establishes the mean-square stability of *stand-alone* adaptive agents for small step-sizes (Lemma II.1), and provides expressions for their MSD and ER metrics (Lemma II.2). We now examine two situations involving a multitude of similar agents. In the first scenario, each agent senses data and analyzes it independently of the other agents. We refer to this mode of operation as non-cooperative processing. In the second scenario, the agents

transmit the collected data for processing at a fusion center. We refer to this mode of operation as centralized or batch processing. We motivate the discussion by considering first the case of mean-square-error costs. Subsequently, we extend the results to more general costs.

A. Non-Cooperative Mean-Square-Error (MSE) Processing

Thus, consider *separate* agents, labeled $k = 1, 2, \dots, N$. Each agent, k , receives streaming data

$$\{\mathbf{d}_k(i), \mathbf{u}_{k,i}; i \geq 0\} \quad (48a)$$

where we are using the subscript k to index the data at agent k . We assume that the data at each agent satisfies the same statistical properties as in Example II.3, and the same linear regression model (13a) with a common w^o albeit with noise $\mathbf{v}_k(i)$. We denote the statistical moments of the data at agent k by

$$R_{u,k} = \mathbb{E} \mathbf{u}_{k,i}^\top \mathbf{u}_{k,i} > 0, \quad \sigma_{v,k}^2 = \mathbb{E} \mathbf{v}_k^2(i) \quad (48b)$$

We further assume in this motivating example that the $R_{u,k}$ are uniform across the agents so that

$$R_{u,k} \equiv R_u, \quad k = 1, 2, \dots, N \quad (48c)$$

In this way, the cost $J_k(w) = \mathbb{E}(\mathbf{d}_k(i) - \mathbf{u}_{k,i}w)^2$, which is associated with agent k , will satisfy a condition similar to (9) with the corresponding parameters $\{\nu, \delta\}$ given by (cf. (10a)):

$$\nu = 2\lambda_{\min}(R_u), \quad \delta = 2\lambda_{\max}(R_u) \quad (49a)$$

Now, if each agent runs the asynchronous LMS learning rule (44) to estimate w^o on its own then, according to (45a), each agent k will attain an individual MSD level that is given by

$$\text{MSD}_{\text{ncop},k}^{\text{asyn}} \approx \mu_x M \sigma_{v,k}^2, \quad k = 1, 2, \dots, N \quad (49b)$$

where we are further assuming that the parameter μ_x is uniform across the agents to enable a meaningful comparison. Moreover, according to (28b), agent k will converge towards this level at a rate dictated by:

$$\alpha_{\text{ncop},k}^{\text{asyn}} \approx 1 - 4\mu_x \lambda_{\min}(R_u) \quad (49c)$$

The subscript “ncop” is used in (49b) and (49c) to indicate that these expressions are for the non-cooperative mode of operation. It is seen from (49b) that agents with noisier data (i.e., larger $\sigma_{v,k}^2$) will perform worse and have larger MSD levels than agents with cleaner data. In other words, whenever adaptive agents act individually, the quality of their solution will be as good as the quality of their noisy data. This is a sensible conclusion. We are going to show in later sections that cooperation among the agents, whereby agents share information with their neighbors, can help enhance their individual performance levels.

B. Centralized Mean-Square-Error (MSE) Processing

Let us now contrast the above non-cooperative solution with a centralized implementation whereby, at every iteration i , the

N agents transmit their raw data $\{\mathbf{d}_k(i), \mathbf{u}_{k,i}\}$ to a fusion center for processing. In a *synchronous* environment, once the fusion center receives the raw data, it can run a standard stochastic-gradient update of the form:

$$\mathbf{w}_i = \mathbf{w}_{i-1} + \mu \left[\frac{1}{N} \sum_{k=1}^N 2\mathbf{u}_{k,i}^\top (\mathbf{d}_k(i) - \mathbf{u}_{k,i}\mathbf{w}_{i-1}) \right] \quad (50a)$$

where μ is the constant step-size, and the term multiplying μ can be seen to correspond to a sample average of several approximate gradient vectors. The analysis in [1], [2] showed that the MSD performance that results from this implementation is given by (using expression (65a) with $H_k = 2R_u$ and $R_{s,k} = 4\sigma_{v,k}^2 R_u$):

$$\text{MSD}_{\text{cent}}^{\text{sync}} \approx \frac{\mu M}{N} \left(\frac{1}{N} \sum_{k=1}^N \sigma_{v,k}^2 \right) \quad (50b)$$

Moreover, using expression (65b) given further ahead, this centralized solution will converge towards the above MSD level at the same rate as the non-cooperative solution:

$$\alpha_{\text{cent}}^{\text{sync}} \approx 1 - 4\mu \lambda_{\min}(R_u) \quad (50c)$$

In an asynchronous environment, there are now *several* random events that can interfere with the operation of the fusion center. Let us consider initially one particular random event that corresponds to the situation in which the fusion center may or may not update at any particular iteration (e.g., due to some power-saving strategy). In a manner similar to (21), we may represent this scenario by writing:

$$\mathbf{w}_i = \mathbf{w}_{i-1} + \mu(i) \left[\frac{1}{N} \sum_{k=1}^N 2\mathbf{u}_{k,i}^\top (\mathbf{d}_k(i) - \mathbf{u}_{k,i}\mathbf{w}_{i-1}) \right] \quad (51a)$$

with a random step-size process, $\mu(i)$; this process is again assumed to satisfy the conditions under Assumption II.3. The analysis in the sequel will show that the MSD performance that results from this implementation is given by:

$$\text{MSD}_{\text{cent}}^{\text{asyn},1} \approx \frac{\mu_x M}{N} \left(\frac{1}{N} \sum_{k=1}^N \sigma_{v,k}^2 \right) \quad (51b)$$

where we are adding the superscript “1” to indicate that this is a preliminary result pertaining to the particular asynchronous implementation (51a). We will be generalizing this result very soon, at which point we will drop the superscript “1”. Likewise, using expression (65b) given further ahead, this version of the asynchronous centralized solution converges towards the above MSD level at the same rate as the non-cooperative solution (49c) and the synchronous version (50c):

$$\alpha_{\text{cent}}^{\text{asyn},1} \approx 1 - 4\mu_x \lambda_{\min}(R_u) \quad (51c)$$

Observe from (50b) and (51b) that the MSD level attained by the centralized solution is proportional to $1/N$ times the *average* noise power across all agents. This scaled average noise power can be larger than some of the individual noise variances and smaller than the remaining noise variances. For

example, consider a situation with $N = 2$ agents, $\sigma_{v,2}^2 = 5\sigma_v^2$ and $\sigma_{v,1}^2 = \sigma_v^2$. Then,

$$\frac{1}{N} \left(\frac{1}{N} \sum_{k=1}^N \sigma_{v,k}^2 \right) = \frac{1}{2} \left(\frac{1}{2} (\sigma_{v,1}^2 + \sigma_{v,2}^2) \right) = 1.5\sigma_v^2 \quad (52)$$

which is larger than $\sigma_{v,1}^2$ and smaller than $\sigma_{v,2}^2$. In this case, the asynchronous centralized solution performs better than the non-cooperative asynchronous agent 2 (i.e., leads to a smaller MSD) but worse than the non-cooperative asynchronous agent 1. This example shows that it does not generally hold that centralized stochastic-gradient implementations outperform all individual non-cooperative agents [77].

Now, more generally, observe that in the synchronous batch processing case (50a), as well as in the asynchronous implementation (51a), the data collected from the various agents are equally aggregated with a weighting factor equal to $1/N$. We can incorporate a second type of random events besides the random variations in the step-size parameter. This second source of uncertainty involves the possibility of failure (or weakening) in the links connecting the individual agents to the fusion center (e.g., due to fading or outage, or perhaps due to the agents themselves deciding to enter into a sleep mode according to some power-saving policy). We can capture these possibilities by extending formulation (51a) in the following manner:

$$\mathbf{w}_i = \mathbf{w}_{i-1} + \boldsymbol{\mu}(i) \left[\sum_{k=1}^N 2\pi_k(i) \mathbf{u}_{k,i}^\top (\mathbf{d}_k(i) - \mathbf{u}_{k,i} \mathbf{w}_{i-1}) \right] \quad (53a)$$

where $\boldsymbol{\mu}(i)$ continues to be a random step-size process, but now the coefficients $\{\pi_k(i); k = 1, 2, \dots, N\}$ are new random fusion coefficients that satisfy:

$$\sum_{k=1}^N \pi_k(i) = 1, \quad \pi_k(i) \geq 0 \quad (53b)$$

for every $i \geq 0$.

Assumption III.1 (Conditions on random fusion coefficients). *It is assumed that the random fusion coefficients $\{\pi_k(i)\}$ are independent of any other random variable in the learning algorithm, and that $\pi_k(i)$ and $\pi_\ell(j)$ are independent of each other for any $i \neq j$ and $k \neq \ell$. This latter condition means that the process $\{\pi_k(i)\}$ is assumed to be independent over both time and space. The mean and co-variance(s) of each $\pi_k(i)$ are denoted by*

$$\bar{\pi}_k \triangleq \mathbb{E} \pi_k(i) \quad (54a)$$

$$c_{\pi,k\ell} \triangleq \mathbb{E} (\pi_k(i) - \bar{\pi}_k)(\pi_\ell(i) - \bar{\pi}_\ell) \quad (54b)$$

for all $k, \ell = 1, 2, \dots, N$ and all $i \geq 0$. When $k = \ell$, the scalar $c_{\pi,kk}$ corresponds to the variance of $\pi_k(i)$ and, therefore, we shall also use the alternative notation $\sigma_{\pi,k}^2$ for this case:

$$\sigma_{\pi,k}^2 \triangleq c_{\pi,kk} \geq 0 \quad (54c)$$

□

It is straightforward to verify that the first and second-order

moments of the coefficients $\{\pi_k(i)\}$ satisfy:

$$\bar{\pi}_k \geq 0, \quad \sum_{k=1}^N \bar{\pi}_k = 1 \quad (55a)$$

$$\sum_{k=1}^N c_{\pi,k\ell} = 0, \quad \text{for any } \ell \quad (55b)$$

$$\sum_{\ell=1}^N c_{\pi,k\ell} = 0, \quad \text{for any } k \quad (55c)$$

Note that the earlier asynchronous implementation (51a) is a special case of the more general formulation (53a) when the mean of the random fusion coefficients are chosen as $\bar{\pi}_k = 1/N$ and their variances are set to zero, $\sigma_{\pi,k}^2 = 0$. In order to enable a fair and meaningful comparison among the three centralized implementations (50a), (51a), and (53a), we shall assume that the means of the fusion coefficients in the latest implementation are set to $\bar{\pi}_k = 1/N$ (results for arbitrary mean values are listed in Example III.1 further ahead and appear in [5]). We will show in the sequel that for this choice of $\bar{\pi}_k$, the MSD performance of implementation (53a) is given by:

$$\text{MSD}_{\text{cent}}^{\text{asyn}} \approx \frac{\mu_x M}{N} \left[\frac{1}{N} \sum_{k=1}^N (1 + N^2 \sigma_{\pi,k}^2) \sigma_{v,k}^2 \right] \quad (56a)$$

Moreover, using expression (63b) given further ahead, this asynchronous centralized solution will converge towards the above MSD level at the same rate as the non-cooperative solution (49c) and the synchronous version (50c):

$$\alpha_{\text{cent}}^{\text{asyn}} \approx 1 - 4\mu_x \lambda_{\min}(R_u) \quad (56b)$$

C. Stochastic-Gradient Centralized Solution

The previous two sections focused on mean-square-error costs. We now extend the conclusions to more general costs. Thus, consider a collection of N agents, each with an individual convex cost function, $J_k(w)$. The objective is to determine the unique minimizer w^o of the aggregate cost:

$$J^{\text{glob}}(w) \triangleq \sum_{k=1}^N J_k(w) \quad (57)$$

It is now the above aggregate cost, $J^{\text{glob}}(w)$, that will be required to satisfy the conditions of Assumptions II.1 and II.4 relative to some parameters $\{\nu_c, \delta_c, \tau_c\}$, with the subscript “c” used to indicate that these factors correspond to the centralized implementation. Under these conditions, the cost $J^{\text{glob}}(w)$ will have a unique minimizer, which we continue to denote by w^o . We will not be requiring each individual cost, $J_k(w)$, to be strongly convex. It is sufficient for at least one of these costs to be strongly convex while the remaining costs can be convex; this condition ensures the strong convexity of $J^{\text{glob}}(w)$. Although some individual costs may not have a unique minimizer, we require in this exposition that w^o is one of their minima so that all individual costs share a minimum at location w^o ; the treatment in [1], [2], [74] considers the

more general case in which the minimizers of the individual costs $\{J_k(w)\}$ can be different and need not contain a common location, w^o .

There are many centralized solutions that can be used to determine the unique minimizer w^o of (57), with some solution techniques being more powerful than other techniques. Nevertheless, we shall focus on centralized implementations of the stochastic-gradient type. The reason we consider the *same* class of stochastic gradient algorithms for non-cooperative, centralized, and distributed solutions in this chapter is to enable a *meaningful* comparison among the various implementations. Thus, we first consider a *synchronous* centralized strategy of the following form:

$$w_i = w_{i-1} - \mu \left[\frac{1}{N} \sum_{k=1}^N \widehat{\nabla_{w^\top} J_k}(w_{i-1}) \right], \quad i \geq 0 \quad (58a)$$

with a constant step-size, μ . When the fusion center employs random step-sizes, the above solution is replaced by:

$$w_i = w_{i-1} - \mu(i) \left[\frac{1}{N} \sum_{k=1}^N \widehat{\nabla_{w^\top} J_k}(w_{i-1}) \right], \quad i \geq 0 \quad (58b)$$

where the process $\mu(i)$ now satisfies Assumption II.3. More generally, the asynchronous implementation can employ random fusion coefficients as well such as:

$$w_i = w_{i-1} - \mu(i) \sum_{k=1}^N \pi_k(i) \widehat{\nabla_{w^\top} J_k}(w_{i-1}), \quad i \geq 0 \quad (58c)$$

where the coefficients $\{\pi_k(i)\}$ satisfy Assumption III.1 with means

$$\bar{\pi}_k = 1/N \quad (58d)$$

D. Performance of Centralized Solution

To examine the performance of the asynchronous implementation (58c)–(58d), we proceed in two steps. First, we identify the gradient noise that is present in the recursion; it is equal to the difference between the true gradient vector for the global cost, $J^{\text{glob}}(w)$, defined by (57) and its approximation. Second, we argue that (58c) has a form similar to the single-agent stochastic-gradient algorithm (21) and, therefore, invoke earlier results to write down performance metrics for the centralized solution (58c)–(58d).

We start by introducing the individual gradient noise processes:

$$s_{k,i}(w_{i-1}) \triangleq \widehat{\nabla_{w^\top} J_k}(w_{i-1}) - \nabla_{w^\top} J_k(w_{i-1}) \quad (59a)$$

for $k = 1, 2, \dots, N$. We assume that these noises satisfy conditions similar to Assumption II.2 with parameters $\{\beta_k^2, \sigma_{s,k}^2, R_{s,k}\}$, i.e.,

$$R_{s,k} \triangleq \lim_{i \rightarrow \infty} \mathbb{E} [s_{k,i}(w^o) s_{k,i}^\top(w^o) | \mathcal{F}_{i-1}] \quad (59b)$$

and

$$\mathbb{E} [s_{k,i}(w_{i-1}) | \mathcal{F}_{i-1}] = 0 \quad (59c)$$

$$\mathbb{E} [\|s_{k,i}(w_{i-1})\|^2 | \mathcal{F}_{i-1}] \leq \beta_k^2 \|\tilde{w}_{i-1}\|^2 + \sigma_{s,k}^2 \quad (59d)$$

Additionally, we assume that the gradient noise components across the agents are uncorrelated with each other:

$$\mathbb{E} [s_{k,i}(w_{i-1}) s_{\ell,i}^\top(w_{i-1}) | \mathcal{F}_{i-1}] = 0, \quad \text{all } k \neq \ell \quad (59e)$$

Using these gradient noise terms, it is straightforward to verify that recursion (58c) can be rewritten as:

$$w_i = w_{i-1} - \frac{\mu(i)}{N} \left[s_i(w_{i-1}) + \sum_{k=1}^N \nabla_{w^\top} J_k(w_{i-1}) \right] \quad (60)$$

where $s_i(w_{i-1})$ denotes the overall gradient noise; its expression is given by

$$s_i(w_{i-1}) = \sum_{k=1}^N \left[N \pi_k(i) \widehat{\nabla_{w^\top} J_k}(w_{i-1}) - \nabla_{w^\top} J_k(w_{i-1}) \right] \quad (61)$$

Since iteration (60) has the form of a stochastic gradient recursion with random update similar to (21), we can infer its mean-square-error behavior from Lemmas II.1 and II.2 if the noise process $s_i(w_{i-1})$ can be shown to satisfy conditions similar to Assumption II.2 with some parameters $\{\beta_c^2, \sigma_s^2\}$. Indeed, starting from (61), some algebra will show that

$$\mathbb{E} [s_i(w_{i-1}) | \mathcal{F}_{i-1}] = 0 \quad (62a)$$

$$\mathbb{E} [\|s_i(w_{i-1})\|^2 | \mathcal{F}_{i-1}] \leq \beta_c^2 \|\tilde{w}_{i-1}\|^2 + \sigma_s^2 \quad (62b)$$

where

$$\beta_c^2 \triangleq \sum_{k=1}^N [\beta_k^2 + N^2 \sigma_{\pi,k}^2 (\beta_k^2 + \delta_c^2)] \quad (62c)$$

$$\sigma_s^2 \triangleq \sum_{k=1}^N (1 + N^2 \sigma_{\pi,k}^2) \sigma_{s,k}^2 \quad (62d)$$

The following result now follows from Lemmas II.1 and II.2 [2], [5].

Lemma III.1 (Convergence of centralized solution). *Assume the aggregate cost (57) satisfies the conditions under Assumption II.1 for some parameters $0 < \nu_c \leq \delta_c$. Assume also that the individual gradient noise processes defined by (59a) satisfy the conditions under Assumption II.2 for some parameters $\{\beta_k^2, \sigma_{s,k}^2, R_{s,k}\}$, in addition to the orthogonality condition (59e). Let $\mu_o = 2\nu_c/(\delta_c^2 + \beta_c^2)$. For any $\mu_x/N < \mu_o$, the iterates generated by the asynchronous centralized solution (58c)–(58d) satisfy:*

$$\mathbb{E} \|\tilde{w}_i\|^2 \leq \alpha \mathbb{E} \|\tilde{w}_{i-1}\|^2 + \sigma_s^2 (\bar{\mu}^2 + \sigma_\mu^2) / N^2 \quad (63a)$$

where the scalar α satisfies $0 \leq \alpha < 1$ and is given by

$$\alpha_{\text{cent}}^{\text{asyn}} = 1 - 2\nu_c (\mu_x/N) + (\delta_c^2 + \beta_c^2) (\bar{\mu}^2 + \sigma_\mu^2) / N^2 \quad (63b)$$

It follows from (63a) that for sufficiently small step-size parameter $\mu_x \ll 1$:

$$\limsup_{i \rightarrow \infty} \mathbb{E} \|\tilde{w}_i\|^2 = O(\mu_x) \quad (63c)$$

Moreover, under smoothness conditions similar to (35a) for $J^{\text{glob}}(w)$ for some parameter $\tau_c \geq 0$, and similar to (35b) for the individual

gradient noise covariance matrices, it holds for small μ_x that:

$$\text{MSD}_{\text{cent}}^{\text{asyn}} = \frac{\mu_x}{2N} \text{Tr} \left[\left(\sum_{k=1}^N H_k \right)^{-1} \left(\sum_{k=1}^N (1 + N^2 \sigma_{\pi,k}^2) R_{s,k} \right) \right] + O(\mu_x^{3/2}) \quad (64)$$

where $H_k = \nabla_w^2 J_k(w^o)$. \square

We can recover from the expressions in the lemma, performance results for the particular asynchronous implementation described earlier by (58b) by setting $\sigma_{\pi,k}^2 = 0$ so that

$$\text{MSD}_{\text{cent}}^{\text{asyn},1} \approx \frac{\mu_x}{2N} \text{Tr} \left[\left(\sum_{k=1}^N H_k \right)^{-1} \left(\sum_{k=1}^N R_{s,k} \right) \right] \quad (65a)$$

$$\alpha_{\text{cent}}^{\text{asyn},1} \approx 1 - 2\nu_c \mu_x / N \quad (65b)$$

It is seen from (64) and (65a) that when the fusion center operates under the more general asynchronous policy (58c), the additional randomness in the fusion coefficients $\{\pi_k(i)\}$ degrades the MSD performance relative to (65a) due to the presence of the factors $1 + N^2 \sigma_{\pi,k}^2 > 1$, i.e., to first-order in μ_x we have:

$$\text{MSD}_{\text{cent}}^{\text{asyn},1} < \text{MSD}_{\text{cent}}^{\text{asyn}} \quad (66a)$$

In comparison, the added randomness due to the $\{\pi_k(i)\}$ does not have significant impact on the convergence rate since it is straightforward to see that, to first order in the step-size parameter μ_x :

$$\alpha_{\text{cent}}^{\text{asyn},1} \approx \alpha_{\text{cent}}^{\text{asyn}} \quad (66b)$$

Likewise, setting $\bar{\mu}_k = \mu$ and $\sigma_{\pi,k}^2 = 0$, the performance of the synchronous centralized solution (58a) is obtained as a special case of the results in the lemma:

$$\text{MSD}_{\text{cent}}^{\text{sync}} \approx \frac{\mu}{2N} \text{Tr} \left[\left(\sum_{k=1}^N H_k \right)^{-1} \left(\sum_{k=1}^N R_{s,k} \right) \right] \quad (67a)$$

$$\alpha_{\text{cent}}^{\text{sync}} \approx 1 - 2\nu_c \mu / N \quad (67b)$$

These expressions agree with the performance results presented in [1], [2].

Example III.1 (Case of general mean-values). Although not used in this chapter, we remark in passing that if the mean values $\{\bar{\pi}_k\}$ are not fixed at $1/N$, as was required by (58d), then the MSD performance of the asynchronous centralized solution (58c) will instead be given by

$$\text{MSD}_{\text{cent}}^{\text{asyn}} = \frac{\mu_x N}{2} \text{Tr} \left[\left(\sum_{k=1}^N H_k \right)^{-1} \left(\sum_{k=1}^N (\bar{\pi}_k^2 + \sigma_{\pi,k}^2) R_{s,k} \right) \right] + O(\mu_x^{3/2}) \quad (68)$$

◆

E. Comparison with Non-Cooperative Processing

We can now compare the performance of the asynchronous centralized solution (58c) against the performance of non-cooperative processing when agents act independently of each

other and run the recursion:

$$\boxed{w_{k,i} = w_{k,i-1} - \mu(i) \widehat{\nabla_{w^T} J_k(w_{k,i-1})}, \quad i \geq 0} \quad (69a)$$

This comparison is *meaningful* when all agents share the same unique minimizer so that we can compare how well the individual agents are able to recover the same w^o as the centralized solution. For this reason, we re-introduce the requirement that all individual costs $\{J_k(w)\}$ are ν -strongly convex with a uniform parameter ν . Since $J^{\text{glob}}(w)$ is the aggregate sum of the individual costs, then we can set the lower bound ν_c for the Hessian of $J^{\text{glob}}(w)$ at $\nu_c = N\nu$. From expressions (28b) and (65b) we then conclude that, for a sufficiently small μ_x , the convergence rates of the asynchronous non-cooperative solution (69a) and the asynchronous centralized solution with random update (58b) will be similar:

$$\begin{aligned} \alpha_{\text{cent}}^{\text{asyn}} &\approx 1 - 2\nu_c (\mu_x / N) \\ &= 1 - 2\nu \mu_x \\ &\approx \alpha_{\text{ncop},k}^{\text{asyn}} \end{aligned} \quad (69b)$$

Moreover, we observe from (41a) that the average MSD level across N non-cooperative asynchronous agents is given by

$$\text{MSD}_{\text{ncop,av}}^{\text{asyn}} \approx \frac{\mu_x}{2N} \text{Tr} \left[\sum_{k=1}^N H_k^{-1} R_{s,k} \right] \quad (69c)$$

so that comparing with (65a), some simple algebra allows us to conclude that, for small step-sizes and to first order in μ_x :

$$\text{MSD}_{\text{cent}}^{\text{asyn},1} \leq \text{MSD}_{\text{ncop,av}}^{\text{asyn}} \quad (69d)$$

That is, while the asynchronous centralized solution (58b) with random updates need not outperform every individual non-cooperative agent in general, its performance outperforms the average performance across all non-cooperative agents.

The next example illustrates this result by considering the scenario where all agents have the same Hessian matrices at $w = w^o$, namely, $H_k \equiv H$ for $k = 1, 2, \dots, N$. This situation occurs, for example, when the individual costs are identical across the agents, say, $J_k(w) \equiv J(w)$, as is common in machine learning applications. This situation also occurs for the mean-square-error costs we considered earlier in this section when the regression covariance matrices, $\{R_{u,k}\}$, are uniform across all agents, i.e., $R_{u,k} \equiv R_u$ for $k = 1, 2, \dots, N$. In these cases with uniform Hessian matrices H_k , the example below establishes that the asynchronous centralized solution (58b) with random updates improves over the average MSD performance of the non-cooperative solution (69a) by a factor of N .

Example III.2 (N -fold improvement in MSD performance). Consider a collection of N agents whose individual cost functions, $J_k(w)$, are ν -strongly convex and are minimized at the same location $w = w^o$. The costs are also assumed to have identical Hessian matrices at $w = w^o$, i.e., $H_k \equiv H$. Then, using (65a), the MSD of the asynchronous centralized implementation (58b) with random updates is given by:

$$\text{MSD}_{\text{cent}}^{\text{asyn},1} \approx \frac{1}{N} \left(\frac{\mu_x}{2N} \sum_{k=1}^N \text{Tr}(H^{-1} R_{s,k}) \right) \approx \frac{1}{N} \text{MSD}_{\text{ncop,av}}^{\text{asyn}} \quad (70)$$

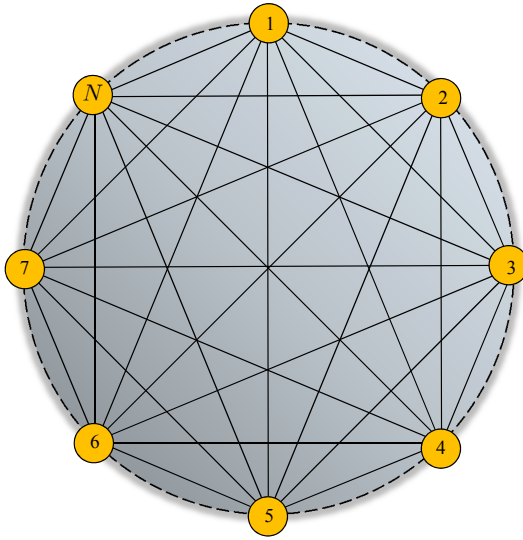


Fig. 3. Example of a fully-connected network, where each agent can access information from all other agents.

Example III.3 (Random fusion can degrade MSD performance). Although the convergence rates of the asynchronous centralized solution (58c) and the non-cooperative solution (69a) agree to first-order in μ_x , the relation between their MSD values is indefinite (contrary to (69d)), as illustrated by the following example.

Consider the same setting of Example III.2 and assume further that the variances of the random fusion coefficients are uniform, i.e., $\sigma_{\pi,k}^2 \equiv \sigma_\pi^2$. Then, using (64), the MSD of the asynchronous centralized implementation (58c) is given by

$$\begin{aligned} \text{MSD}_{\text{cent}}^{\text{asyn}} &\approx \frac{1 + N^2 \sigma_\pi^2}{N} \left(\frac{\mu_x}{2N} \sum_{k=1}^N \text{Tr}(H^{-1} R_{s,k}) \right) \\ &= \left(\frac{1 + N^2 \sigma_\pi^2}{N} \right) \text{MSD}_{\text{ncop,av}}^{\text{asyn}} \end{aligned} \quad (71)$$

Therefore, to first-order in μ_x , we find that

$$\begin{cases} \text{MSD}_{\text{cent}}^{\text{asyn}} \leq \text{MSD}_{\text{ncop,av}}^{\text{asyn}}, & \text{if } \sigma_\pi^2 \leq \frac{N-1}{N^2} \\ \text{MSD}_{\text{cent}}^{\text{asyn}} > \text{MSD}_{\text{ncop,av}}^{\text{asyn}}, & \text{if } \sigma_\pi^2 > \frac{N-1}{N^2} \end{cases} \quad (72)$$

In other words, if the variance of the random fusion coefficients is large enough, then the centralized solution will generally have degraded performance relative to the non-cooperative solution (which is an expected result).

Example III.4 (Fully-connected networks). In preparation for the discussion on networked agents, it is useful to describe one extreme situation where a collection of N agents are fully connected to each other — see Fig. 3. In this case, each agent is able to access the data from all other agents and, therefore, each individual agent can run a synchronous or asynchronous centralized implementation, say, one of the same form as (58b):

$$\mathbf{w}_{k,i} = \mathbf{w}_{k,i-1} - \boldsymbol{\mu}(i) \left[\frac{1}{N} \sum_{\ell=1}^N \widehat{\nabla_{\mathbf{w}^\top} J_\ell}(\mathbf{w}_{k,i-1}) \right], \quad i \geq 0 \quad (73)$$

When this happens, each agent will attain the same performance level as that of the asynchronous centralized solution (58b). Two observations are in place. First, note from (73) that the information that agent k is receiving from all other agents is their gradient

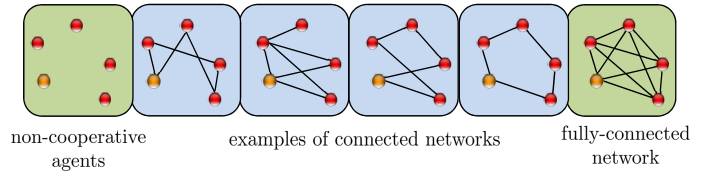


Fig. 4. Examples of connected networks, with the left-most panel representing a collection of non-cooperative agents.

vector approximations. Obviously, other pieces of information could be shared among the agents, such as their iterates $\{\mathbf{w}_{\ell,i-1}\}$. Second, note that the right-most term multiplying $\boldsymbol{\mu}(i)$ in (73) corresponds to a convex combination of the approximate gradients from the various agents, with the combination coefficients being uniform and all equal to $1/N$. In general, there is no need for these combination weights to be identical. Even more importantly, agents do not need to have access to information from all other agents in the network. We are going to see in the sequel that interactions with a limited number of neighbors is sufficient for the agents to attain performance that is comparable to that of the centralized solution.

Figure 4 shows a simple selection of connected topologies for five agents. The leftmost panel corresponds to the non-cooperative case and the rightmost panel corresponds to the fully-connected case. The panels in between illustrate some other topologies. In the coming sections, we are going to present results that allow us to answer useful questions about such networked agents such as: (a) Which topology has best performance in terms of mean-square error and convergence rate? (b) Given a connected topology, can it be made to approach the performance of the centralized solution? (c) Which aspects of the topology influence performance? (d) Which aspects of the combination weights (policy) influence performance? (e) Can different topologies deliver similar performance levels? (f) Is cooperation always beneficial? and (g) If the individual agents are able to solve the inference task individually in a stable manner, does it follow that the connected network will remain stable regardless of the topology and regardless of the cooperation strategy?

IV. SYNCHRONOUS MULTI-AGENT ADAPTATION AND LEARNING

In this section, we describe distributed strategies of the consensus (e.g., [15], [19], [23], [78]–[89]) and diffusion (e.g., [1], [2], [14], [15], [32], [48], [90], [91]) types. These strategies rely solely on localized interactions among neighboring agents, and they can be used to seek the minimizer of (57). We first describe the network model.

A. Strongly-Connected Networks

Figure 5 shows an example of a network consisting of N connected agents, labeled $k = 1, 2, \dots, N$. Following the presentation from [2], [15], the network is represented by a graph consisting of N vertices (representing the agents) and a set of edges connecting the agents to each other. An edge that connects an agent to itself is called a self-loop. The neighborhood of an agent k is denoted by \mathcal{N}_k and it consists of all agents that are connected to k by an edge. Any two neighboring agents k and ℓ have the ability to share information over the edge connecting them.

We assume an undirected graph so that if agent k is a neighbor of agent ℓ , then agent ℓ is also a neighbor of agent k . We assign a pair of nonnegative scaling weights, $\{a_{k\ell}, a_{\ell k}\}$,

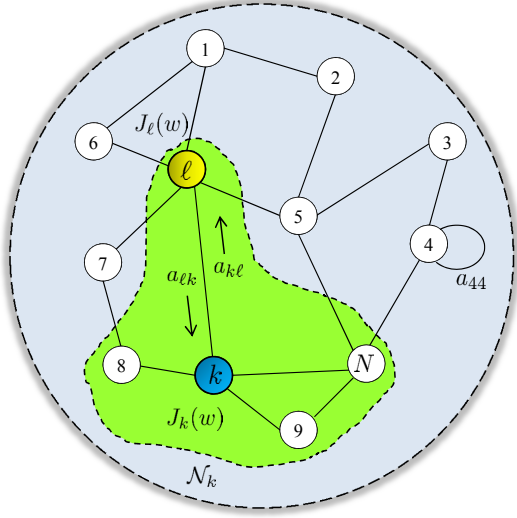


Fig. 5. Agents that are linked by edges can share information. The neighborhood of agent k is marked by the highlighted area.

to the edge connecting k and ℓ . The scalar $a_{\ell k}$ is used by agent k to scale the data it receives from agent ℓ ; this scaling can be interpreted as a measure of the confidence level that agent k assigns to its interaction with agent ℓ . Likewise, $a_{k\ell}$ is used by agent ℓ to scale the data it receives from agent k . The weights $\{a_{k\ell}, a_{\ell k}\}$ can be different so that the exchange of information between the neighboring agents $\{k, \ell\}$ need not be symmetrical. One or both weights can also be zero.

A network is said to be connected if paths with nonzero scaling weights can be found linking any two distinct agents in *both* directions, either directly when they are neighbors or by passing through intermediate agents when they are not neighbors. In this way, information can flow in both directions between any two agents in the network, although the forward path from an agent k to some other agent ℓ need not be the same as the backward path from ℓ to k . A strongly-connected network is a connected network with at least one non-trivial self-loop, meaning that $a_{kk} > 0$ for some agent k .

The strong connectivity of a network translates into a useful property on the combination weights. Assume we collect the coefficients $\{a_{\ell k}\}$ into an $N \times N$ matrix $A = [a_{\ell k}]$, such that the entries on the k -th column of A contain the coefficients used by agent k to scale data arriving from its neighbors $\ell \in \mathcal{N}_k$; we set $a_{\ell k} = 0$ if $\ell \notin \mathcal{N}_k$. We refer to A as the *combination matrix* or *policy*. It turns out that combination matrices that correspond to strongly-connected networks are *primitive* — an $N \times N$ matrix A with nonnegative entries is said to be primitive if there exists some finite integer $n_o > 0$ such that all entries of A^{n_o} are strictly positive [2], [15], [92].

B. Distributed Optimization

Network cooperation can be exploited to solve adaptation, learning, and optimization problems in a decentralized manner in response to streaming data. To explain how cooperation can be achieved, we start by associating with each agent k a twice-differentiable cost function $J_k(w) : \mathbb{R}^{M \times 1} \mapsto \mathbb{R}$.

The objective of the network of agents is to seek the unique minimizer of the aggregate cost function, $J^{\text{glob}}(w)$, defined by (57). Now, however, we seek a *distributed* (as opposed to a centralized) solution. In a distributed implementation, each agent k can only rely on its own data and on data from its neighbors.

We continue to assume that $J^{\text{glob}}(w)$ satisfies the conditions of Assumptions II.1 and II.4 with parameters $\{\nu_d, \delta_d, \tau_d\}$, with the subscript “d” now used to indicate that these parameters are related to the distributed implementation. Under these conditions, the cost $J^{\text{glob}}(w)$ will have a unique minimizer, which we continue to denote by w^o . For simplicity of presentation, we will also assume in the remainder of this chapter that the individual costs $J_k(w)$ are strongly convex as well. These costs can be distinct across the agents or they can all be identical, i.e., $J_k(w) \equiv J(w)$ for $k = 1, 2, \dots, N$; in the latter situation, the problem of minimizing (57) would correspond to the case in which the agents work together to optimize the same cost function. If we let w_k^o denote the minimizer of $J_k(w)$, we continue to assume for this exposition that each $J_k(w)$ is also minimized at w^o :

$$w_k^o \equiv w^o, \quad k = 1, 2, \dots, N \quad (74)$$

The case where the individual costs are only convex and need not be strongly convex is discussed in [1], [2], [49]; most of the results and conclusions continue to hold but the derivations become more technical. Likewise, the case in which the individual costs need not share minimizers is discussed in these references. There are nevertheless important situations in practice where the minimizers of individual costs coincide with each other. For instance, examples abound where agents need to work cooperatively to attain a common objective such as tracking a target, locating a food source, or evading a predator (see, e.g., [14], [93], [94]). This scenario is also common in machine learning problems [38], [39], [95]–[98] when data samples at the various agents are generated by a common distribution parameterized by some vector, w^o . One such situation is illustrated in the next example.

Example IV.1 (Mean-square-error (MSE) networks). Consider the same setting of Example II.3 except that we now have N agents observing streaming data $\{\mathbf{d}_k(i), \mathbf{u}_{k,i}\}$ that satisfy the regression model (13a) with regression covariance matrices $R_{u,k} = \mathbb{E} \mathbf{u}_{k,i}^T \mathbf{u}_{k,i} > 0$ and with the same unknown w^o , i.e.,

$$\mathbf{d}_k(i) = \mathbf{u}_{k,i} w^o + \mathbf{v}_k(i) \quad (75a)$$

The individual mean-square-error costs are defined by

$$J_k(w) = \mathbb{E} (\mathbf{d}_k(i) - \mathbf{u}_{k,i} w)^2 \quad (75b)$$

and are strongly convex in this case, with the minimizer of each $J_k(w)$ occurring at

$$w_k^o \triangleq R_{u,k}^{-1} r_{du,k}, \quad k = 1, 2, \dots, N \quad (75c)$$

If we multiply both sides of (75a) by $\mathbf{u}_{k,i}^T$ from the left, and take expectations, we find that w^o satisfies

$$r_{du,k} = R_{u,k} w^o \quad (76)$$

This relation shows that the unknown w^o from (75a) satisfies the same expression as w_k^o in (75c), for any $k = 1, 2, \dots, N$, so that we must have $w^o = w_k^o$. Therefore, this example amounts to a situation

where all costs $\{J_k(w)\}$ attain their minima at the same location, w^o .

We shall use the network model of this example to illustrate other results in the chapter. For ease of reference, we shall refer to strongly-connected networks with agents receiving data according to model (75a) and seeking to estimate w^o by adopting the mean-square-error costs $J_k(w)$ defined above, as *mean-square-error (MSE) networks*. We assume for these networks that the measurement noise process $v_k(i)$ is temporally white and independent over space so that

$$\mathbb{E} v_k(i) v_\ell(j) = \sigma_{v,k}^2 \delta_{k,\ell} \delta_{i,j} \quad (77)$$

in terms of the Kronecker delta $\delta_{k,\ell}$. Likewise, we assume that the regression data $u_{k,i}$ is temporally white and independent over space so that

$$\mathbb{E} u_{k,i}^\top u_{\ell,j} = R_{u,k} \delta_{k,\ell} \delta_{i,j} \quad (78)$$

Moreover, the measurement noise $v_k(i)$ and the regression data $u_{\ell,j}$ are independent of each other for all k, ℓ, i, j . These statistical conditions help facilitate the analysis of such networks.

In the next subsections we list *synchronous* distributed algorithms of the consensus and diffusion types for the optimization of $J^{\text{glob}}(w)$. We only list the algorithms here; for motivation and justifications, the reader may refer to the treatments in [1], [2].

C. Synchronous Consensus Strategy

Let $w_{k,i}$ denote the iterate that is available at agent k at iteration i ; this iterate serves as the estimate for w^o . The consensus iteration at each agent k is described by the following construction:

$$w_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} w_{\ell,i-1} - \mu_k \widehat{\nabla_{w^\top} J_k}(w_{k,i-1}) \quad (79)$$

where the $\{\mu_k\}$ are individual step-size parameters, and where the combination coefficients $\{a_{\ell k}\}$ that appear in (79) are nonnegative scalars that are required to satisfy the following conditions for each agent $k = 1, 2, \dots, N$:

$$a_{\ell k} \geq 0, \quad \sum_{\ell=1}^N a_{\ell k} = 1, \quad a_{\ell k} = 0 \text{ if } \ell \notin \mathcal{N}_k \quad (80a)$$

Condition (80a) implies that the combination matrix $A = [a_{\ell k}]$ satisfies

$$A^\top \mathbf{1} = \mathbf{1} \quad (80b)$$

where $\mathbf{1}$ denotes the vector with all entries equal to one. We say that A is left-stochastic. One useful property of left-stochastic matrices is that their spectral radius is equal to one [15], [92], [99]–[101]:

$$\rho(A) = 1 \quad (80c)$$

An equivalent representation that is useful for later analysis is to rewrite the consensus iteration (79) as shown in the following listing, where the intermediate iterate that results from the neighborhood combination is denoted by $\psi_{k,i-1}$. Observe that the gradient vector in the consensus implementation (81) is evaluated at $w_{k,i-1}$ and not $\psi_{k,i-1}$.

Consensus strategy for distributed adaptation

for each time instant $i \geq 0$:

each agent $k = 1, 2, \dots, N$ performs the update:

$$\begin{cases} \psi_{k,i-1} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} w_{\ell,i-1} \\ w_{k,i} = \psi_{k,i-1} - \mu_k \widehat{\nabla_{w^\top} J_k}(w_{k,i-1}) \end{cases} \quad (81)$$

end

We remark that one way to motivate the consensus update (79) is to start from the non-cooperative step (69a) and replace the first iterate $w_{k,i-1}$ by the convex combination used in (79).

Example IV.2 (Consensus LMS networks). For the MSE network of Example IV.1, the consensus strategy reduces to:

$$\begin{cases} \psi_{k,i-1} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} w_{\ell,i-1} \\ w_{k,i} = \psi_{k,i-1} + 2\mu_k u_{k,i}^\top [d_k(i) - u_{k,i} w_{k,i-1}] \end{cases} \quad (82)$$

D. Synchronous Diffusion Strategies

There is an inherent asymmetry in the consensus construction. Observe from the computation of $w_{k,i}$ in (81) that the update starts from $\psi_{k,i-1}$ and corrects it by the approximate gradient vector evaluated at $w_{k,i-1}$ (and not at $\psi_{k,i-1}$). This *asymmetry* will be shown later, e.g., in Example VII.4, to be problematic when the consensus strategy is used for adaptation and learning over networks. This is because the asymmetry can cause an unstable growth in the state of the network [32]. Diffusion strategies remove the asymmetry problem.

Combine-then-Adapt (CTA) Diffusion. In the CTA formulation of the diffusion strategy, the *same* iterate $\psi_{k,i-1}$ is used to compute $w_{k,i}$, thus leading to description (83a) where the gradient vector is evaluated at $\psi_{k,i-1}$ as well. The reason for the name ‘‘Combine-then-Adapt’’ is that the first step in (83a) involves a combination step, while the second step involves an adaptation step. The reason for the qualification ‘‘diffusion’’ is that the use of $\psi_{k,i-1}$ to evaluate the gradient vector allows information to diffuse more thoroughly through the network. This is because information is not only being diffused through the aggregation of the neighborhood iterates, but also through the evaluation of the gradient vector at the aggregate state value.

Diffusion strategy for distributed adaptation (CTA)

for each time instant $i \geq 0$:

each agent $k = 1, 2, \dots, N$ performs the update:

$$\begin{cases} \psi_{k,i-1} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} w_{\ell,i-1} \\ w_{k,i} = \psi_{k,i-1} - \mu_k \widehat{\nabla_{w^\top} J_k}(\psi_{k,i-1}) \end{cases} \quad (83a)$$

end

Adapt-then-Combine (ATC) Diffusion. A similar implementa-

tion can be obtained by switching the order of the combination and adaptation steps in (83a), as shown in the listing (83b). The structure of the CTA and ATC strategies are fundamentally identical: the difference lies in which variable we choose to correspond to the updated iterate $\mathbf{w}_{k,i}$. In ATC, we choose the result of the *combination* step to be $\mathbf{w}_{k,i}$, whereas in CTA we choose the result of the *adaptation* step to be $\mathbf{w}_{k,i}$.

Diffusion strategy for distributed adaptation (ATC)

for each time instant $i \geq 0$:

each agent $k = 1, 2, \dots, N$ performs the update:

$$\begin{cases} \boldsymbol{\psi}_{k,i} = \mathbf{w}_{k,i-1} - \mu_k \widehat{\nabla_{\mathbf{w}^\top} J_k}(\mathbf{w}_{k,i-1}) \\ \mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \boldsymbol{\psi}_{\ell,i} \end{cases} \quad (83b)$$

end

Example IV.3 (Diffusion LMS networks). For the MSE network of Example IV.1, the ATC and CTA diffusion strategies reduce to:

$$\begin{cases} \boldsymbol{\psi}_{k,i-1} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \mathbf{w}_{\ell,i-1} & (\text{CTA diffusion}) \\ \mathbf{w}_{k,i} = \boldsymbol{\psi}_{k,i-1} + 2\mu_k \mathbf{u}_{k,i}^\top [\mathbf{d}_k(i) - \mathbf{u}_{k,i} \boldsymbol{\psi}_{k,i-1}] \end{cases} \quad (84a)$$

and

$$\begin{cases} \boldsymbol{\psi}_{k,i} = \mathbf{w}_{k,i-1} + 2\mu_k \mathbf{u}_{k,i}^\top [\mathbf{d}_k(i) - \mathbf{u}_{k,i} \mathbf{w}_{k,i-1}] \\ \mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \boldsymbol{\psi}_{\ell,i} & (\text{ATC diffusion}) \end{cases} \quad (84b)$$

◆

Example IV.4 (Diffusion logistic regression). We revisit the pattern classification problem from Example II.2, where we consider a collection of N networked agents cooperating with each other to solve the logistic regression problem. Each agent receives streaming data $\{\boldsymbol{\gamma}_k(i), \mathbf{h}_{k,i}\}$, where the variable $\boldsymbol{\gamma}_k(i)$ assumes the values ± 1 and designates the class that the feature vector $\mathbf{h}_{k,i}$ belongs to. The objective is to use the training data to determine the vector \mathbf{w}^o that minimizes the cost

$$J_k(\mathbf{w}) \triangleq \frac{\rho}{2} \|\mathbf{w}\|^2 + \mathbb{E} \left\{ \ln \left[1 + e^{-\boldsymbol{\gamma}_k(i) \mathbf{h}_{k,i}^\top \mathbf{w}} \right] \right\} \quad (85)$$

under the assumption of joint wide-sense stationarity over the random data. It is straightforward to verify that the ATC diffusion strategy (83b) reduces to the following form in this case:

$$\begin{cases} \boldsymbol{\psi}_{k,i} = (1 - \rho\mu_k) \mathbf{w}_{k,i-1} + \mu_k \left(\frac{\boldsymbol{\gamma}_k(i)}{1 + e^{\boldsymbol{\gamma}_k(i) \mathbf{h}_{k,i}^\top \mathbf{w}_{k,i-1}}} \right) \mathbf{h}_{k,i} \\ \mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \boldsymbol{\psi}_{\ell,i} \end{cases} \quad (86)$$

◆

Diffusion Strategies with Enlarged Cooperation. Other forms of diffusion cooperation are possible by allowing for enlarged exchange of information among the agents, such as exchanging gradient vector approximations *in addition* to the iterates. For example, the following ATC form (similarly for CTA) employs an additional set of combination coefficients $\{c_{\ell k}\}$ to aggregate gradient information [15], [48], [91]:

$$\begin{cases} \boldsymbol{\psi}_{k,i} = \mathbf{w}_{k,i-1} - \mu_k \sum_{\ell \in \mathcal{N}_k} c_{\ell k} \widehat{\nabla_{\mathbf{w}^\top} J_\ell}(\mathbf{w}_{k,i-1}) \\ \mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \boldsymbol{\psi}_{\ell,i} \end{cases} \quad (\text{ATC diffusion}) \quad (87a)$$

where the $\{c_{\ell k}\}$ are nonnegative scalars that satisfy:

$$c_{\ell k} \geq 0, \quad \sum_{k=1}^N c_{\ell k} = 1, \quad \text{and} \quad c_{\ell k} = 0 \text{ if } \ell \notin \mathcal{N}_k \quad (87b)$$

If we collect the entries $\{c_{\ell k}\}$ into an $N \times N$ matrix C , so that the ℓ -th row of C is formed of $\{c_{\ell k}, k = 1, 2, \dots, N\}$, then $C\mathbf{1} = \mathbf{1}$ and we say that C is a *right-stochastic* matrix. Observe that (87a) is equivalent to associating with each agent k the weighted neighborhood cost function:

$$J'_k(\mathbf{w}) \triangleq \sum_{\ell \in \mathcal{N}_k} c_{\ell k} J_\ell(\mathbf{w}) \quad (88)$$

Our discussion in the sequel focuses on the case $C = I_N$.

E. Discussion and Related Literature

As remarked in [1], [2], [15], there has been extensive work on consensus techniques in the literature, starting with the foundational results by [102], [103], which were of a different nature and did not respond to streaming data arriving continuously at the agents, as is the case, for instance, with the continuous arrival of data $\{\mathbf{d}_k(i), \mathbf{u}_{k,i}\}$ in Examples IV.2 and IV.3.. The original consensus formulation deals instead with the problem of computing averages over graphs. This can be explained as follows [19], [79], [102], [103]. Consider a collection of (scalar or vector) measurements denoted by $\{w_\ell, \ell = 1, 2, \dots, N\}$ available at the vertices of a connected graph with N agents. The objective is to devise a distributed algorithm that enables every agent to determine the *average* value:

$$\bar{w} \triangleq \frac{1}{N} \sum_{k=1}^N w_k \quad (89)$$

by interacting solely with its neighbors. When this occurs, we say that the agents have reached consensus (or agreement) about \bar{w} . We select an $N \times N$ *doubly-stochastic* combination matrix $A = [a_{\ell k}]$; a doubly-stochastic matrix is one that has nonnegative elements and satisfies

$$A^\top \mathbf{1} = \mathbf{1}, \quad A \mathbf{1} = \mathbf{1} \quad (90)$$

We assume the second largest-magnitude eigenvalue of A satisfies

$$|\lambda_2(A)| < 1 \quad (91)$$

Using the combination coefficients $\{a_{\ell k}\}$, each agent k then iterates *repeatedly* on the data of its neighbors:

$$\mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \mathbf{w}_{\ell,i-1}, \quad i \geq 0, \quad k = 1, 2, \dots, N \quad (92)$$

starting from the boundary conditions $\mathbf{w}_{\ell,-1} = \mathbf{w}_\ell$ for all $\ell \in \mathcal{N}_k$. The superscript i continues to denote the iteration index. Every agent k in the network performs the same calculation, which amounts to combining repeatedly, and in

a convex manner, the state values of its neighbors. It can then be shown that (see [102], [103] and [15, App.E]):

$$\lim_{i \rightarrow \infty} w_{k,i} = \bar{w}, \quad k = 1, 2, \dots, N \quad (93)$$

In this way, through the localized iterative process (92), the agents are able to converge to the global average value, \bar{w} .

Motivated by this elegant result, several works in the literature (e.g., [19], [23], [43], [78], [80]–[82], [104]–[110]) proposed useful extensions of the original consensus construction (92) to minimize aggregate costs of the form (57) or to solve distributed estimation problems of the least-squares or Kalman filtering type. Some of the earlier extensions involved the use of *two* separate time-scales: one faster time-scale for performing multiple consensus iterations similar to (92) over the states of the neighbors, and a second slower time-scale for performing gradient vector updates or for updating the estimators by using the result of the consensus iterations (e.g., [78], [82], [105]–[109]). An example of a two-time scale implementation would be an algorithm of the following form:

$$\begin{cases} w_{\ell,i-1}^{(-1)} \leftarrow w_{\ell,i-1}, \text{ for all agents } \ell \text{ at iteration } i-1 \\ \text{for } n = 0, 1, 2, \dots, J-1 \text{ iterate:} \\ \quad w_{k,i-1}^{(n)} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} w_{\ell,i-1}^{(n-1)}, \text{ for all } k = 1, 2, \dots, N \\ \text{end} \\ w_{k,i} = w_{k,i-1}^{(J-1)} - \mu_k \widehat{\nabla_{w^*} J}_k(w_{k,i-1}) \end{cases} \quad (94)$$

If we compare the last equation in (94) with (81), we observe that the variable $w_{k,i-1}^{(J-1)}$ that is used in (94) to obtain $w_{k,i}$ is the result of J repeated applications of a consensus operation of the form (92) on the iterates $\{w_{\ell,i-1}\}$. The purpose of these repeated calculations is to approximate well the average of the iterates in the neighborhood of agent k . These J repeated averaging operations need to be completed before the availability of the gradient information for the last update step in (94). In other words, the J averaging operations need to be performed at a faster rate than the last step in (94). Such two time-scale implementations are a hindrance for real-time adaptation from streaming data. The separate time-scales turn out to be unnecessary and this fact was one of the motivations for the introduction of single time-scale diffusion strategies in [90], [111]–[118].

Building upon a useful procedure for distributed optimization from [19, Eq. (2.1)] and [43, Eq. (7.1)], more recent works proposed single time-scale implementations for consensus strategies as well by using an implementation similar to (79) — see, e.g., [89, Eq. (3)], [81, Eq. (3)], [83, Eq. (19)], and [23, Eq.(9)]. These references, however, generally employ decaying step-sizes, $\mu_k(i) \rightarrow 0$, to ensure that the iterates $\{w_{k,i}\}$ across all agents will converge almost-surely to the same value (thus, reaching agreement or consensus), namely, they employ recursions of the form:

$$w_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} w_{\ell,i-1} - \mu_k(i) \widehat{\nabla_{w^*} J}_k(w_{k,i-1}) \quad (95)$$

or variations thereof, such as replacing $\mu_k(i)$ by some time-

variant gain matrix sequence, say, $K_{k,i}$:

$$w_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} w_{\ell,i-1} - K_{k,i} \cdot \widehat{\nabla_{w^*} J}_k(w_{k,i-1}) \quad (96)$$

As noted before, when diminishing step-sizes are used, adaptation is turned off over time, which is prejudicial for learning purposes. For this reason, we are instead setting the step-sizes to constant values in (79) in order to endow the consensus iteration with continuous adaptation and learning abilities (and to enhance the convergence rate). It turns out that some care is needed for consensus implementations when constant step-sizes are used. The main reason is that, as already explained in [1] and [2][Sec. 10.6], and as alluded to earlier, instability can occur in consensus networks due to an inherent asymmetry in the dynamics of the consensus iteration — see also Example VII.4 further ahead.

A second main reason for the introduction of cooperative strategies of the diffusion type (83a) and (83b) has been to show that single time-scale distributed learning from *streaming* data is possible, and that this objective can be achieved under *constant* step-size adaptation in a stable manner [3], [74], [90], [91], [111]–[114], [116] — see also [2][Chs. 9–11]; the diffusion strategies further allow A to be left-stochastic and permit larger modes of cooperation than doubly-stochastic policies. The CTA diffusion strategy (84a) was first introduced for mean-square-error estimation problems in [90], [111]–[113]. The ATC diffusion structure (84b), with adaptation preceding combination, appeared in the work [115] on adaptive distributed least-squares schemes and also in the works [91], [116]–[118] on distributed mean-square-error and state-space estimation methods. The CTA structure (83a) with an iteration dependent step-size that decays to zero, $\mu(i) \rightarrow 0$, was employed in [21], [119], [120] to solve distributed optimization problems that require all agents to reach agreement. The ATC form (83b), also with an iteration dependent sequence $\mu(i)$ that decays to zero, was employed in [121], [122] to ensure almost-sure convergence and agreement among agents.

V. ASYNCHRONOUS MULTI-AGENT ADAPTATION AND LEARNING

There are various ways by which asynchronous events can be introduced into the operation of a distributed strategy. Without loss in generality, we illustrate the model for asynchronous operation by describing it for the ATC diffusion strategy (83b); similar constructions apply to CTA diffusion (83a) and consensus (81).

A. Asynchronous Model

In a first instance, we model the step-size parameters as random variables and replace (83b) by:

$$\begin{cases} \psi_{k,i} = w_{k,i-1} - \mu_k(i) \widehat{\nabla_{w^*} J}_k(w_{k,i-1}) \\ w_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \psi_{\ell,i} \end{cases} \quad (97)$$

In this model, the neighborhoods and the network topology remain fixed and only the $\mu_k(i)$ assume random values. The step-sizes can vary across the agents and, therefore, their

means and variances become agent-dependent. Moreover, the step-sizes across agents can be correlated with each other. We therefore denote the first and second-order moments of the step-size parameters by:

$$\bar{\mu}_k \triangleq \mathbb{E} \mu_k(i) \quad (98a)$$

$$\sigma_{\mu,k}^2 \triangleq \mathbb{E} (\mu_k(i) - \bar{\mu}_k)^2 \quad (98b)$$

$$c_{\mu,k\ell} \triangleq \mathbb{E} (\mu_k(i) - \bar{\mu}_k)(\mu_\ell(i) - \bar{\mu}_\ell) \quad (98c)$$

When $\ell = k$, the scalar $c_{\mu,kk}$ coincides with the variance of $\mu_k(i)$, i.e., $c_{\mu,kk} = \sigma_{\mu,k}^2 \geq 0$. On the other hand, if the step-sizes across the agents happen to be uncorrelated, then $c_{\mu,k\ell} = 0$ for $k \neq \ell$.

More broadly, we can allow for random variations in the neighborhoods (and, hence, in the network structure), and random variations in the combination coefficients as well. We capture this more general asynchronous implementation by writing:

$$\begin{cases} \psi_{k,i} = \mathbf{w}_{k,i-1} - \mu_k(i) \widehat{\nabla_{\mathbf{w}^\top} J_k}(\mathbf{w}_{k,i-1}) \\ \mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_{k,i}} a_{\ell k}(i) \psi_{\ell,i} \end{cases} \quad (99)$$

where the combination coefficients $\{a_{\ell k}(i)\}$ are now *random* and, moreover, the symbol $\mathcal{N}_{k,i}$ denotes the *randomly-changing* neighborhood of agent k at time i . These neighborhoods become random because the random variations in the combination coefficients can turn links on and off depending on the values of the $\{a_{\ell k}(i)\}$. We continue to require the combination coefficients $\{a_{\ell k}(i)\}$ to satisfy the same structural constraint as given before by (80a), i.e.,

$$\sum_{\ell \in \mathcal{N}_{k,i}} a_{\ell k}(i) = 1, \text{ and } \begin{cases} a_{\ell k}(i) > 0, & \text{if } \ell \in \mathcal{N}_{k,i} \\ a_{\ell k}(i) = 0, & \text{otherwise} \end{cases} \quad (100)$$

Since these coefficients are now random, we denote their first and second-order moments by:

$$\bar{a}_{\ell k} \triangleq \mathbb{E} a_{\ell k}(i) \quad (101a)$$

$$\sigma_{a,\ell k}^2 \triangleq \mathbb{E} (a_{\ell k}(i) - \bar{a}_{\ell k})^2 \quad (101b)$$

$$c_{a,\ell k,nm} \triangleq \mathbb{E} (a_{\ell k}(i) - \bar{a}_{\ell k})(a_{nm}(i) - \bar{a}_{nm}) \quad (101c)$$

When $\ell = n$ and $k = m$, the scalar $c_{a,\ell k,nm}$ coincides with the variance of $a_{\ell k}(i)$, i.e., $c_{a,\ell k,nm} = \sigma_{a,\ell k}^2 \geq 0$.

Example V.1 (Asynchronous diffusion LMS networks). For the MSE network of Example IV.1, the ATC diffusion strategy (97) with random update reduces to

$$\begin{cases} \psi_{k,i} = \mathbf{w}_{k,i-1} + 2\mu_k(i) \mathbf{u}_{k,i}^\top [\mathbf{d}_k(i) - \mathbf{u}_{k,i} \mathbf{w}_{k,i-1}] \\ \mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \psi_{\ell,i} \quad (\text{diffusion with random updates}) \end{cases} \quad (102a)$$

whereas the *asynchronous* ATC diffusion strategy (99) reduces to:

$$\begin{cases} \psi_{k,i} = \mathbf{w}_{k,i-1} + 2\mu_k(i) \mathbf{u}_{k,i}^\top [\mathbf{d}_k(i) - \mathbf{u}_{k,i} \mathbf{w}_{k,i-1}] \\ \mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_{k,i}} a_{\ell k}(i) \psi_{\ell,i} \quad (\text{asynchronous diffusion}) \end{cases} \quad (102b)$$

Obviously, implementation (102a) is a special case of the asynchronous update (102b) when the variances of the random combination coefficients $\{a_{\ell k}(i)\}$ are set to zero.

Example V.2 (Asynchronous consensus LMS networks). Similarly, for the same MSE network of Example IV.1, the asynchronous consensus strategy is given by

$$\begin{cases} \psi_{k,i-1} = \sum_{\ell \in \mathcal{N}_{k,i}} a_{\ell k}(i) \mathbf{w}_{\ell,i-1} \\ \mathbf{w}_{k,i} = \psi_{k,i-1} + 2\mu_k(i) \mathbf{u}_{k,i}^\top [\mathbf{d}_k(i) - \mathbf{u}_{k,i} \mathbf{w}_{k,i-1}] \end{cases} \quad (103)$$

B. Mean Graph

We refer to the topology that corresponds to the average combination coefficients $\{\bar{a}_{\ell k}\}$ as the *mean graph*, which is fixed over time. For each agent k , the neighborhood defined by the mean graph is denoted by \mathcal{N}_k . It is straightforward to verify that the mean combination coefficients $\bar{a}_{\ell k}$ satisfy the following constraints over the mean graph (compare with (80a) and (100)):

$$\sum_{\ell \in \mathcal{N}_k} \bar{a}_{\ell k} = 1, \text{ and } \begin{cases} \bar{a}_{\ell k} > 0, & \text{if } \ell \in \mathcal{N}_k \\ \bar{a}_{\ell k} = 0, & \text{otherwise} \end{cases} \quad (104)$$

One example of a random network with two equally probable realizations and its mean graph is shown in Fig. 6 [3]. The letter ω is used to index the sample space of the random matrix \mathbf{A}_i .

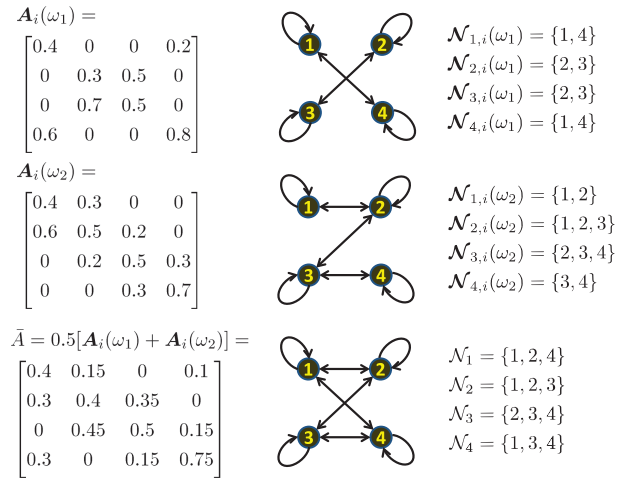


Fig. 6. The first two rows show two equally probable realizations with the respective neighborhoods. The last row shows the resulting mean graph.

There is one useful result that relates the random neighborhoods $\{\mathcal{N}_{k,i}\}$ from (99) to the neighborhoods $\{\mathcal{N}_k\}$ from

the mean graph. It is not difficult to verify that \mathcal{N}_k is equal to the *union* of all possible realizations for the random neighborhoods $\mathcal{N}_{k,i}$; this property is already illustrated by the example of Fig. 6:

$$\mathcal{N}_k = \bigcup_{\omega \in \Omega} \mathcal{N}_{k,i}(\omega) \quad (105)$$

for any k , where Ω denotes the sample space for $\mathcal{N}_{k,i}$.

C. Random Combination Policy

The first and second-order moments of the combination coefficients will play an important role in characterizing the stability and mean-square-error performance of the asynchronous network (99). We collect these moments in matrix form as follows. We first group the combination coefficients into a matrix:

$$\mathbf{A}_i \triangleq [\mathbf{a}_{\ell k}(i)]_{\ell,k=1}^N \quad (N \times N) \quad (106a)$$

The sequence $\{\mathbf{A}_i, i \geq 0\}$ represents a stochastic process consisting of left-stochastic random matrices whose entries satisfy the conditions in (100) at every time i . We subsequently introduce the mean and Kronecker-covariance matrix of \mathbf{A}_i and assume these quantities are constant over time; we denote them by the $N \times N$ matrix \bar{A} and the $N^2 \times N^2$ matrix C_A , respectively:

$$\bar{A} \triangleq \mathbb{E} \mathbf{A}_i = [\bar{a}_{\ell k}]_{\ell,k=1}^N \quad (107a)$$

$$C_A \triangleq \mathbb{E}[(\mathbf{A}_i - \bar{A}) \otimes (\mathbf{A}_i - \bar{A})] \quad (107b)$$

The matrix C_A is *not* a conventional covariance matrix and is not necessarily Hermitian. The reason for its introduction is because it captures the correlations of each entry of \mathbf{A}_i with all other entries in \mathbf{A}_i . For example, for a network with $N = 2$ agents, the entries of \bar{A} and C_A will be given by:

$$\bar{A} = \begin{bmatrix} \bar{a}_{11} & \bar{a}_{12} \\ \bar{a}_{21} & \bar{a}_{22} \end{bmatrix} \quad (108a)$$

$$C_A = \begin{bmatrix} c_{a,11,11} & c_{a,11,12} & c_{a,12,11} & c_{a,12,12} \\ c_{a,11,21} & c_{a,11,22} & c_{a,12,21} & c_{a,12,22} \\ c_{a,21,11} & c_{a,21,12} & c_{a,22,11} & c_{a,22,12} \\ c_{a,21,21} & c_{a,21,22} & c_{a,22,21} & c_{a,22,22} \end{bmatrix} \quad (108b)$$

We thus see that the (ℓ, k) -th block of C_A contains the covariance coefficients of $\mathbf{a}_{\ell k}$ with all other entries of \mathbf{A}_i . One useful property of the matrices $\{\bar{A}, C_A\}$ so defined is that their elements are nonnegative and the following matrices are left-stochastic:

$$(\bar{A})^\top \mathbf{1}_N = \mathbf{1}_N, \quad (\bar{A} \otimes \bar{A} + C_A)^\top \mathbf{1}_{N^2} = \mathbf{1}_{N^2} \quad (109)$$

Assumption V.1 (Asynchronous network model). *It is assumed that the random processes $\{\mu_k(i), \mathbf{a}_{\ell m}(j)\}$ are independent of each other for all k, ℓ, m, i , and j . They are also independent of any other random variable in the learning algorithm.* \square

The asynchronous network model described in this section covers many situations of practical interest. For example, we

can choose the sample space for each step-size $\mu_k(i)$ to be the binary choice $\{0, \mu\}$ to model random “on-off” behavior at each agent k for the purpose of saving power, waiting for data, or even due to random agent failures. Similarly, we can choose the sample space for each combination coefficient $\mathbf{a}_{\ell k}(i)$, $\ell \in \mathcal{N}_k \setminus \{k\}$, to be $\{0, a_{\ell k}\}$ to model a random “on-off” status for the link from agent ℓ to agent k at time i for the purpose of either saving communication cost or due to random link failures. Note that the convex constraint (100) can always be satisfied by adjusting the value of $\mathbf{a}_{\ell k}(i)$ according to the realizations of $\{\mathbf{a}_{\ell k}(i); \ell \in \mathcal{N}_{k,i} \setminus \{k\}\}$.

Example V.3 (The spatially-uncorrelated model). A useful special case of the asynchronous network model of this section is the spatially-uncorrelated model. In this case, at each iteration i , the random step-sizes $\{\mu_k(i); k = 1, 2, \dots, N\}$ are uncorrelated with each other across the network, and the random combination coefficients $\{\mathbf{a}_{\ell k}(i); \ell \neq k, k = 1, 2, \dots, N\}$ are also uncorrelated with each other across the network. Then, it can be verified that the covariances $\{c_{\mu, k\ell}\}$ in (98c) and $\{c_{a, \ell k, nm}\}$ in (101c) will be fully determined by the variances $\{\sigma_{\mu, k}^2, \sigma_{a, \ell k}^2\}$:

$$c_{\mu, kk} = \sigma_{\mu, k}^2, \quad c_{\mu, k\ell} = 0, \quad k \neq \ell \quad (110a)$$

and

$$c_{a, \ell k, nm} = \begin{cases} \sigma_{a, \ell k}^2, & \text{if } k = m, \ell = n, \ell \in \mathcal{N}_k \setminus \{k\} \\ -\sigma_{a, \ell k}^2, & \text{if } k = m = n, \ell \in \mathcal{N}_k \setminus \{k\} \\ -\sigma_{a, nk}^2, & \text{if } k = m = \ell, n \in \mathcal{N}_k \setminus \{k\} \\ \sum_{j \in \mathcal{N}_k \setminus \{k\}} \sigma_{a, jk}^2, & \text{if } k = m = \ell = n \\ 0, & \text{otherwise} \end{cases} \quad (110b)$$

D. Perron Vectors

Now that we introduced the network model, we can move on to examine the effect of network cooperation on performance. Some interesting patterns of behavior arise when agents cooperate to solve a global optimization problem in a distributed manner from streaming data [1], [2]. For example, one interesting result established in [3], [4], [49], [74], [123] is that the effect of the network topology on performance is captured by the Perron vector of the combination policy. This vector turns out to summarize the influence of the topology on performance so much so that different topologies with similar Perron vectors will end up delivering similar performance. We explained the role of Perron vectors in the context of synchronous adaptation and learning in [1], [2]. Here we focus on asynchronous networks. In this case, two Perron vectors will be needed since the randomness in the combination policy is now represented by two moment matrices, \bar{A} and C_A . In the synchronous case, only one Perron vector was necessary since the combination policy was fixed and described by a matrix A . Although we are focusing on the asynchronous case in the sequel, we will be able to recover results for synchronous networks as special cases.

Let us first recall the definition of Perron vectors for synchronous networks, say, of the form described by (83b) with combination policy A . We assume the network is strongly-connected. In this case, the left-stochastic matrix A will be

primitive. For such primitive matrices, it follows from the Perron-Frobenius Theorem [92] that:

- (a) The matrix A will have a *single* eigenvalue at one.
- (b) All other eigenvalues of A will be strictly inside the unit circle so that $\rho(A) = 1$.
- (c) With proper sign scaling, all entries of the right-eigenvector of A corresponding to the single eigenvalue at one will be *positive*. Let p denote this right-eigenvector with its entries $\{p_k\}$ normalized to add up to one, i.e.,

$$\boxed{Ap = p, \quad \mathbf{1}^\top p = 1, \quad p_k > 0, \quad k = 1, 2, \dots, N} \quad (111)$$

We refer to p as the *Perron* eigenvector of A . It was explained in [1], [2] how the entries of this vector determine the mean-square-error performance and convergence rate of the network; these results will be revisited further ahead when we recover them as special cases of the asynchronous results.

On the other hand, for an asynchronous implementation, the individual realizations of the random combination matrix A_i in (99) need not be primitive. In this context, we will require a form of primitiveness to hold on average as follows.

Definition V.1 (Strongly-connected asynchronous model). We say that an asynchronous model with random combination coefficients $\{\mathbf{a}_{\ell k}(i)\}$ is *strongly-connected* if the Kronecker-covariance matrix given by $\bar{A} \otimes \bar{A} + C_A$ is primitive.

□

Observe that if we set $\bar{A} = A$ and $C_A = 0$, then we recover the condition for strong-connectedness in the synchronous case, namely, that A should be primitive.

Definition V.1 means that the directed graph (digraph) associated with the matrix $\bar{A} \otimes \bar{A} + C_A$ is strongly-connected (e.g., [100, pp. 30,34] and [2]). It is straightforward to check from the definition of C_A in (107b) that

$$\bar{A} \otimes \bar{A} + C_A = \mathbb{E}(\mathbf{A}_i \otimes \mathbf{A}_i) \quad (112)$$

so that the digraph associated with $\bar{A} \otimes \bar{A} + C_A$ is the union of all possible digraphs associated with the realizations of $\mathbf{A}_i \otimes \mathbf{A}_i$ [124, p. 29]. Therefore, as explained in [3]–[5], definition V.1 amounts to an assumption that the *union* of all possible digraphs associated with the realizations of $\mathbf{A}_i \otimes \mathbf{A}_i$ is strongly-connected. As illustrated in Fig. 7, this condition still allows individual digraphs associated with realizations of \mathbf{A}_i to be weakly-connected with or without self-loops or even to be disconnected [4].

It follows from property (109) and Definition V.1, the matrix $\bar{A} \otimes \bar{A} + C_A$ is left-stochastic and primitive. It can be verified that mean matrix \bar{A} is also primitive if $\bar{A} \otimes \bar{A} + C_A$ is primitive (although the converse is not true). Therefore, the matrix \bar{A} is both left-stochastic and primitive. We denote the Perron eigenvector of $\bar{A} \otimes \bar{A} + C_A$ by $p_c \in \mathbb{R}^{N^2 \times 1}$, which satisfies:

$$\boxed{(\bar{A} \otimes \bar{A} + C_A)p_c = p_c, \quad p_c^\top \mathbf{1} = 1} \quad (113a)$$

Likewise, we denote the Perron eigenvector of \bar{A} by $\bar{p} \in$

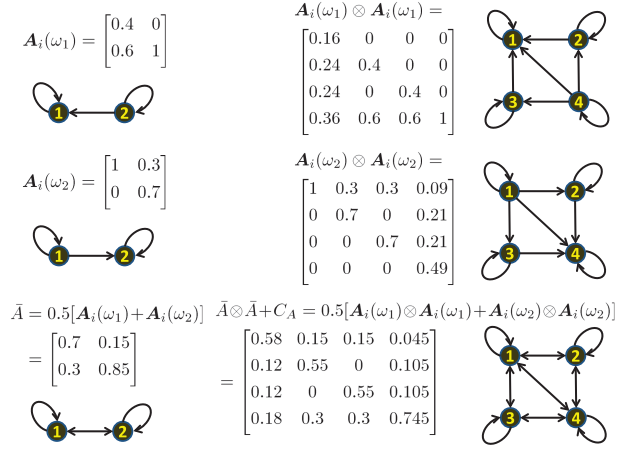


Fig. 7. The combination policy A_i has two equally probable realizations in this example, denoted by $\{A_i(\omega_1), A_i(\omega_2)\}$. Observe that neither of the digraphs $A_i(\omega_1) \otimes A_i(\omega_1)$ or $A_i(\omega_2) \otimes A_i(\omega_2)$ is strongly-connected due to the existence of the source and sink nodes. However, the digraph associated with $\mathbb{E}(A_i \otimes A_i)$, which is the union of the first two digraphs, is strongly-connected, where information can flow in any direction through the network.

$\mathbb{R}^{N \times 1}$, which satisfies:

$$\boxed{\bar{A}\bar{p} = \bar{p}, \quad \bar{p}^\top \mathbf{1} = 1} \quad (113b)$$

Observe that if we set $C_A = 0$ and $\bar{A} = A$, then we recover the synchronous case information, namely, $\bar{p} = p$ and $p_c = p \otimes p$.

The entries of the Perron vectors $\{p_c, \bar{p}\}$ are related to each other, as illustrated by the following explanation (proofs appear in [4][App. VII]). Since the vector p_c is of dimension $N^2 \times 1$, we partition it into N sub-vectors of dimension $N \times 1$ each:

$$p_c \triangleq \text{col}\{p_1, p_2, \dots, p_N\} \quad (114a)$$

where p_k denotes the k -th sub-vector. We construct an $N \times N$ matrix P_c from these sub-vectors:

$$P_c \triangleq [p_1 \ p_2 \ \dots \ p_N] = [p_{c,\ell k}]_{k,\ell=1}^N \quad (114b)$$

We use $p_{\ell k}$ to denote the (ℓ, k) -th element of matrix P_c , which is equal to the ℓ -th element of p_k . It can be verified that the matrix P_c in (114b) is symmetric positive semi-definite and it satisfies

$$\boxed{P_c \mathbf{1} = \bar{p}, \quad P_c^\top = P_c, \quad P_c \geq 0} \quad (114c)$$

where \bar{p} is the Perron eigenvector in (113b). We can further establish the following useful relations:

$$p_{c,\ell k} = p_{c,k\ell}, \quad \sum_{k=1}^N p_{c,\ell k} = \bar{p}_\ell, \quad \sum_{\ell=1}^N p_{c,\ell k} = \bar{p}_k \quad (114d)$$

We can also verify that the matrix difference

$$C_c \triangleq P_c - \bar{p}\bar{p}^\top, \quad C_c^\top = C_c, \quad C_c \geq 0 \quad (114e)$$

is symmetric, positive semi-definite, and satisfies $C_c \mathbf{1} = 0$.

Moreover, it is straightforward to verify that

$$\boxed{c_{c,kk} = p_{c,kk} - \bar{p}_k^2 \geq 0} \quad (115)$$

where $c_{c,kk}$ denotes the (k, k) -th entry in C_c .

VI. ASYNCHRONOUS NETWORK PERFORMANCE

We now comment on the performance of asynchronous networks and compare their metrics against both non-cooperative and centralized strategies.

A. MSD Performance

We denote the MSD performance of the individual agents and the average MSD performance across the network by:

$$\text{MSD}_{\text{dist},k} \triangleq \lim_{i \rightarrow \infty} \mathbb{E} \|\tilde{w}_{k,i}\|^2 \quad (116a)$$

$$\text{MSD}_{\text{dist,av}} \triangleq \frac{1}{N} \sum_{k=1}^N \text{MSD}_{\text{dist},k} \quad (116b)$$

where the error vectors are measured relative to the global optimizer, w^o . We further denote the gradient noise process at the individual agents by

$$\mathbf{s}_{k,i}(w) \triangleq \widehat{\nabla_{w^\top} J_k(w)} - \nabla_{w^\top} J_k(w) \quad (117a)$$

and define

$$H_k \triangleq \nabla_w^2 J_k(w^o) \quad (117b)$$

$$R_{s,k} \triangleq \lim_{i \rightarrow \infty} \mathbb{E} [\mathbf{s}_{k,i}(w^o) \mathbf{s}_{k,i}^\top(w^o) | \mathcal{F}_{i-1}] \quad (117c)$$

where \mathcal{F}_{i-1} now represents the collection of all random events generated by the iterates from across all agents, $\{\mathbf{w}_{k,j}, k = 1, 2, \dots, N\}$, up to time $i - 1$:

$$\mathcal{F}_{i-1} \triangleq \text{filtration}\{\mathbf{w}_{k,-1}, \mathbf{w}_{k,0}, \mathbf{w}_{k,1}, \dots, \mathbf{w}_{k,i-1}, \text{ all } k\} \quad (117d)$$

It is stated later in (149a), under some assumptions on the gradient noise processes (117a), that for strongly-connected asynchronous diffusion networks of the form (97), and for sufficiently small step-sizes μ_x [2], [4]:

$$\text{MSD}_{\text{dist},k}^{\text{asyn}} \approx \text{MSD}_{\text{dist,av}}^{\text{asyn}} \approx \quad (118)$$

$$\frac{1}{2} \text{Tr} \left[\left(\sum_{k=1}^N \bar{\mu}_k \bar{p}_k H_k \right)^{-1} \left(\sum_{k=1}^N (\bar{\mu}_k^2 + \sigma_{\mu,k}^2) p_{c,kk} R_{s,k} \right) \right]$$

The same result holds for asynchronous consensus and CTA diffusion strategies. Observe from (118) the interesting conclusion that the distributed strategy is able to *equalize* the MSD performance across all agents for sufficiently small step-sizes. It is also instructive to compare expression (149a) with (64) and (68) in the centralized case. Observe how cooperation among agents leads to the appearance of the scaling coefficients $\{\bar{p}_k, p_{c,kk}\}$; these factors are determined by \bar{A} and C_A .

Note further that if we set $\bar{\mu}_k = \mu$, $\sigma_{\mu,k}^2 = 0$, $\bar{p}_k = p_k$ and $p_{c,kk} = p_k^2$, then we recover the MSD expression for

synchronous distributed strategies:

$$\text{MSD}_{\text{dist},k}^{\text{sync}} \approx \text{MSD}_{\text{dist,av}}^{\text{sync}} \approx \quad (119)$$

$$\frac{1}{2} \text{Tr} \left[\left(\sum_{k=1}^N \mu_k p_k H_k \right)^{-1} \left(\sum_{k=1}^N \mu_k^2 p_k^2 R_{s,k} \right) \right]$$

This result agrees with expression (62) from [1].

Example VI.1 (MSE networks with random updates). We continue with the setting of Example IV.1, which deals with MSE networks. We assume the first and second-order moments of the random step-sizes are uniform, i.e., $\bar{\mu}_k \equiv \bar{\mu}$ and $c_{\mu,kk} \equiv \sigma_\mu^2$, and also assume uniform regression covariance matrices, i.e., $R_{u,k} \equiv R_u$ for $k = 1, 2, \dots, N$. It follows that $H_k = 2R_u \equiv H$ and $R_{s,k} = 4\sigma_{v,k}^2 R_u$. Substituting into (149a), and assuming a fixed topology with fixed combination coefficients set to $a_{\ell k}$, we conclude that the MSD performance of the diffusion strategy (102a) with random updates is well approximated by:

$$\text{MSD}_{\text{dist},k}^{\text{asyn},1} \approx \text{MSD}_{\text{dist,av}}^{\text{asyn},1} \approx \mu_x M \left(\sum_{k=1}^N p_k^2 \sigma_{v,k}^2 \right) \quad (120a)$$

where

$$\mu_x = \bar{\mu} + \frac{\sigma_\mu^2}{\bar{\mu}} \quad (120b)$$

and p_k is the k -th entry of the Perron vector p defined by (111).

If the combination matrix A happens to be doubly stochastic, then its Perron eigenvector becomes $p = \mathbf{1}/N$. Substituting $p_k = 1/N$ into (120a) gives

$$\text{MSD}_{\text{dist},k}^{\text{asyn},1} \approx \text{MSD}_{\text{dist,av}}^{\text{asyn},1} \approx \frac{\mu_x M}{N} \left(\frac{1}{N} \sum_{k=1}^N \sigma_{v,k}^2 \right) \quad (120c)$$

which agrees with the centralized performance (51b). In other words, the asynchronous diffusion strategy is able to match the performance of the centralized solution for doubly stochastic combination policies, when both implementations employ random updates. Since the centralized solution can improve the average MSD performance over non-cooperative networks, we further conclude that the diffusion strategy can also exceed the average performance of non-cooperative networks. \blacklozenge

Example VI.2 (Asynchronous MSE networks). We continue with the setting of Example VI.1 except that we now employ the asynchronous LMS diffusion network (102b). Its MSD performance is well approximated by:

$$\text{MSD}_{\text{dist},k}^{\text{asyn}} \approx \text{MSD}_{\text{dist,av}}^{\text{asyn}} \approx \mu_x M \left(\sum_{k=1}^N p_{c,kk} \sigma_{v,k}^2 \right) \quad (121a)$$

If the mean combination matrix \bar{A} happens to be doubly stochastic, then its Perron eigenvector becomes $\bar{p} = \mathbf{1}/N$. Substituting $\bar{p}_k = 1/N$ into (120a), and using $p_{c,kk} = \bar{p}_k^2 + c_{c,kk}$, where $c_{c,kk}$ is from (114e), gives

$$\text{MSD}_{\text{dist},k}^{\text{asyn}} \approx \text{MSD}_{\text{dist,av}}^{\text{asyn}} \approx \frac{\mu_x M}{N} \left[\frac{1}{N} \sum_{k=1}^N (1 + N^2 c_{c,kk}) \sigma_{v,k}^2 \right] \quad (121b)$$

It is clear that if $c_{c,kk} = \sigma_{\pi,k}^2$, then the MSD performance in (121b) will agree with the centralized performance (56a). In other words, the distributed diffusion strategy is able to match the performance of the centralized solution. \blacklozenge

VII. NETWORK STABILITY AND PERFORMANCE

In this section, we examine more closely the performance and stability results that were alluded to in the earlier sections. We first examine the consensus and diffusion strategies in a unified manner, and subsequently focus on diffusion strategies due to their enhanced stability properties, as the ensuing discussion will reveal.

A. MSE Networks

We motivate the discussion by presenting first some illustrative examples with MSE networks, which involve quadratic costs. Following the examples, we extend the framework to more general costs.

Example VII.1 (Error dynamics over MSE networks). We consider the MSE network of Example IV.1, which involves quadratic costs with a common minimizer, w^o . The update equations for the non-cooperative, consensus, and diffusion strategies are given by (69a), (82), and (84a)–(84b). We can group these strategies into a single unifying description by considering the following structure in terms of three sets of combination coefficients $\{\mathbf{a}_{o,\ell k}(i), \mathbf{a}_{1,\ell k}(i), \mathbf{a}_{2,\ell k}(i)\}$:

$$\begin{cases} \phi_{k,i-1} = \sum_{\ell \in \mathcal{N}_{k,i}} \mathbf{a}_{1,\ell k}(i) \mathbf{w}_{\ell,i-1} \\ \psi_{k,i} = \sum_{\ell \in \mathcal{N}_{k,i}} \mathbf{a}_{o,\ell k}(i) \phi_{\ell,i-1} + 2\mu_k(i) \mathbf{u}_{k,i}^\top [\mathbf{d}_k(i) - \mathbf{u}_{k,i} \phi_{k,i-1}] \\ \mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_{k,i}} \mathbf{a}_{2,\ell k}(i) \psi_{\ell,i} \end{cases} \quad (122)$$

In (122), the quantities $\{\phi_{k,i-1}, \psi_{k,i}\}$ denote $M \times 1$ intermediate variables, while the nonnegative entries of the $N \times N$ matrices $\mathbf{A}_{o,i} = [\mathbf{a}_{o,\ell k}(i)]$, $\mathbf{A}_{1,i} = [\mathbf{a}_{1,\ell k}(i)]$, and $\mathbf{A}_{2,i} = [\mathbf{a}_{2,\ell k}(i)]$ are assumed to satisfy the same conditions (100). Any of the combination weights $\{\mathbf{a}_{o,\ell k}(i), \mathbf{a}_{1,\ell k}(i), \mathbf{a}_{2,\ell k}(i)\}$ is zero whenever $\ell \notin \mathcal{N}_{k,i}$. Different choices for $\{\mathbf{A}_{o,i}, \mathbf{A}_{1,i}, \mathbf{A}_{2,i}\}$, including random and deterministic choices, correspond to different strategies, as the following examples reveal:

$$\text{non-cooperative: } \mathbf{A}_{1,i} = \mathbf{A}_{o,i} = \mathbf{A}_{2,i} = \mathbf{I}_N \quad (123a)$$

$$\text{consensus: } \mathbf{A}_{o,i} = \mathbf{A}_i, \mathbf{A}_{1,i} = \mathbf{I}_N = \mathbf{A}_{2,i} \quad (123b)$$

$$\text{CTA diffusion: } \mathbf{A}_{1,i} = \mathbf{A}_i, \mathbf{A}_{2,i} = \mathbf{I}_N = \mathbf{A}_{o,i} \quad (123c)$$

$$\text{ATC diffusion: } \mathbf{A}_{2,i} = \mathbf{A}_i, \mathbf{A}_{1,i} = \mathbf{I}_N = \mathbf{A}_{o,i} \quad (123d)$$

where \mathbf{A}_i denotes some generic combination policy satisfying (100). We associate with each agent k the following three errors:

$$\tilde{\mathbf{w}}_{k,i} \triangleq \mathbf{w}^o - \mathbf{w}_{k,i} \quad (124a)$$

$$\tilde{\psi}_{k,i} \triangleq \mathbf{w}^o - \psi_{k,i} \quad (124b)$$

$$\tilde{\phi}_{k,i-1} \triangleq \mathbf{w}^o - \phi_{k,i-1} \quad (124c)$$

which measure the deviations from the global minimizer, w^o . Subtracting w^o from both sides of the equations in (122) we get

$$\begin{cases} \tilde{\phi}_{k,i-1} = \sum_{\ell \in \mathcal{N}_{k,i}} \mathbf{a}_{1,\ell k}(i) \tilde{\mathbf{w}}_{\ell,i-1} \\ \tilde{\psi}_{k,i} = \sum_{\ell \in \mathcal{N}_{k,i}} \mathbf{a}_{o,\ell k}(i) \tilde{\phi}_{\ell,i-1} - 2\mu_k(i) \mathbf{u}_{k,i}^\top \mathbf{u}_{k,i} \tilde{\phi}_{k,i-1} - 2\mu_k(i) \mathbf{u}_{k,i}^\top \mathbf{v}_k(i) \\ \tilde{\mathbf{w}}_{k,i} = \sum_{\ell \in \mathcal{N}_{k,i}} \mathbf{a}_{2,\ell k}(i) \tilde{\psi}_{\ell,i} \end{cases}$$

In a manner similar to (16a), the gradient noise process at each agent k is given by

$$\mathbf{s}_{k,i}(\phi_{k,i-1}) = 2 \left(R_{u,k} - \mathbf{u}_{k,i}^\top \mathbf{u}_{k,i} \right) \tilde{\phi}_{k,i-1} - 2\mathbf{u}_{k,i}^\top \mathbf{v}_k(i) \quad (125a)$$

In order to examine the evolution of the error dynamics across the network, we collect the error vectors from all agents into $N \times 1$ block error vectors (whose individual entries are of size $M \times 1$ each):

$$\tilde{\mathbf{w}}_i \triangleq \begin{bmatrix} \tilde{\mathbf{w}}_{1,i} \\ \tilde{\mathbf{w}}_{2,i} \\ \vdots \\ \tilde{\mathbf{w}}_{N,i} \end{bmatrix}, \quad \tilde{\psi}_i \triangleq \begin{bmatrix} \tilde{\psi}_{1,i} \\ \tilde{\psi}_{2,i} \\ \vdots \\ \tilde{\psi}_{N,i} \end{bmatrix}, \quad \tilde{\phi}_{i-1} \triangleq \begin{bmatrix} \tilde{\phi}_{1,i-1} \\ \tilde{\phi}_{2,i-1} \\ \vdots \\ \tilde{\phi}_{N,i-1} \end{bmatrix} \quad (126a)$$

Motivated by the last term in the second equation in (125a), and by the gradient noise terms (125a), we also introduce the following $N \times 1$ column vectors whose entries are of size $M \times 1$ each:

$$\mathbf{z}_i \triangleq \begin{bmatrix} 2\mathbf{u}_{1,i}^\top \mathbf{v}_1(i) \\ 2\mathbf{u}_{2,i}^\top \mathbf{v}_2(i) \\ \vdots \\ 2\mathbf{u}_{N,i}^\top \mathbf{v}_N(i) \end{bmatrix}, \quad \mathbf{s}_i \triangleq \begin{bmatrix} \mathbf{s}_{1,i}(\phi_{1,i-1}) \\ \mathbf{s}_{2,i}(\phi_{2,i-1}) \\ \vdots \\ \mathbf{s}_{N,i}(\phi_{N,i-1}) \end{bmatrix} \quad (126b)$$

We further introduce the Kronecker products

$$\mathbf{A}_{o,i} \triangleq \mathbf{A}_{o,i} \otimes \mathbf{I}_M, \quad \mathbf{A}_{1,i} \triangleq \mathbf{A}_{1,i} \otimes \mathbf{I}_M, \quad \mathbf{A}_{2,i} \triangleq \mathbf{A}_{2,i} \otimes \mathbf{I}_M \quad (127a)$$

and the following $N \times N$ block diagonal matrices, whose individual entries are of size $M \times M$ each:

$$\mathbf{M}_i \triangleq \text{diag}\{\mu_1(i)\mathbf{I}_M, \mu_2(i)\mathbf{I}_M, \dots, \mu_N(i)\mathbf{I}_M\} \quad (127b)$$

$$\mathbf{R}_i \triangleq \text{diag}\{2\mathbf{u}_{1,i}^\top \mathbf{u}_{1,i}, 2\mathbf{u}_{2,i}^\top \mathbf{u}_{2,i}, \dots, 2\mathbf{u}_{N,i}^\top \mathbf{u}_{N,i}\} \quad (127c)$$

From (125a) we can then easily conclude that the block network variables satisfy the recursions:

$$\begin{cases} \tilde{\phi}_{i-1} &= \mathbf{A}_{1,i}^\top \tilde{\mathbf{w}}_{i-1} \\ \tilde{\psi}_i &= (\mathbf{A}_{o,i}^\top - \mathbf{M}_i \mathbf{R}_i) \tilde{\phi}_{i-1} - \mathbf{M}_i \mathbf{z}_i \\ \tilde{\mathbf{w}}_i &= \mathbf{A}_{2,i}^\top \tilde{\psi}_i \end{cases} \quad (128a)$$

so that the network weight error vector, $\tilde{\mathbf{w}}_i$, evolves according to:

$$\tilde{\mathbf{w}}_i = \mathbf{A}_{2,i}^\top (\mathbf{A}_{o,i}^\top - \mathbf{M}_i \mathbf{R}_i) \mathbf{A}_{1,i}^\top \tilde{\mathbf{w}}_{i-1} - \mathbf{A}_{2,i}^\top \mathbf{M}_i \mathbf{z}_i \quad (128b)$$

For comparison purposes, if each agent operates individually and uses the non-cooperative strategy (69a), then the weight error vector would instead evolve according to the following recursion:

$$\tilde{\mathbf{w}}_i = (\mathbf{I}_{MN} - \mathbf{M}_i \mathbf{R}_i) \tilde{\mathbf{w}}_{i-1} - \mathbf{M}_i \mathbf{z}_i, \quad i \geq 0 \quad (129)$$

where the matrices $\{\mathbf{A}_{o,i}, \mathbf{A}_{1,i}, \mathbf{A}_{2,i}\}$ do not appear any longer, and with a block diagonal coefficient matrix $(\mathbf{I}_{MN} - \mathbf{M}_i \mathbf{R}_i)$. It is also straightforward to verify that recursion (128b) can be equivalently rewritten in the following form in terms of the gradient noise vector, \mathbf{s}_i , defined by (126b):

$$\tilde{\mathbf{w}}_i = \mathbf{B}_i \tilde{\mathbf{w}}_{i-1} + \mathbf{A}_{2,i}^\top \mathbf{M}_i \mathbf{s}_i \quad (130a)$$

where

$$\mathbf{B}_i \triangleq \mathbf{A}_{2,i}^\top (\mathbf{A}_{o,i}^\top - \mathbf{M}_i \mathbf{R}_i) \mathbf{A}_{1,i}^\top \quad (130b)$$

$$\mathbf{R} \triangleq \mathbb{E} \mathbf{R}_i = \text{diag}\{2R_{u,1}, 2R_{u,2}, \dots, 2R_{u,N}\} \quad (130c)$$

◆

Example VII.2 (Mean-error behavior). We continue with the setting of Example VII.1. In mean-square-error analysis, we are interested in examining how the quantities $\mathbb{E} \tilde{\mathbf{w}}_i$ and $\mathbb{E} \|\tilde{\mathbf{w}}_i\|^2$ evolve over time. If we refer back to the data model described in Example IV.1,

where the regression data $\{\mathbf{u}_{k,i}\}$ were assumed to be temporally white and independent over space, then the stochastic matrix \mathcal{R}_i appearing in (128b)–(129) is seen to be statistically independent of $\tilde{\mathbf{w}}_{i-1}$. We further assume that, in the unified formulation, the entries of the combination policies $\{\mathbf{A}_{o,i}, \mathbf{A}_{1,i}, \mathbf{A}_{2,i}\}$ are independent of each other (as well as over time) and of any other variable in the learning algorithm. Therefore, taking expectations of both sides of these recursions, and invoking the fact that $\mathbf{u}_{k,i}$ and $\mathbf{v}_k(i)$ are also independent of each other and have zero means (so that $\mathbb{E} \mathbf{z}_i = 0$), we conclude that the mean-error vectors evolve according to the following recursions:

$$\mathbb{E} \tilde{\mathbf{w}}_i = \bar{\mathcal{B}} (\mathbb{E} \tilde{\mathbf{w}}_{i-1}) \quad (\text{distributed}) \quad (131a)$$

$$\mathbb{E} \tilde{\mathbf{w}}_i = (I_{MN} - \bar{\mathcal{M}}\mathcal{R}) (\mathbb{E} \tilde{\mathbf{w}}_{i-1}) \quad (\text{non-cooperative}) \quad (131b)$$

where

$$\bar{\mathcal{B}} \triangleq \mathbb{E} \mathcal{B}_i = \bar{\mathcal{A}}_2^\top (\bar{\mathcal{A}}_o^\top - \bar{\mathcal{M}}\mathcal{R}) \bar{\mathcal{A}}_1^\top \quad (132a)$$

$$\bar{\mathcal{M}} = \mathbb{E} \mathcal{M}_i = \text{diag}\{\bar{\mu}_1 I_M, \dots, \bar{\mu}_N I_M\} \quad (132b)$$

$$\bar{\mathcal{A}}_o = \mathbb{E} \mathcal{A}_{o,i} \quad (132c)$$

$$\bar{\mathcal{A}}_1 = \mathbb{E} \mathcal{A}_{1,i} \quad (132d)$$

$$\bar{\mathcal{A}}_2 = \mathbb{E} \mathcal{A}_{2,i} \quad (132e)$$

The matrix $\bar{\mathcal{B}}$ controls the dynamics of the mean weight-error vector for the distributed strategies. Observe, in particular, that $\bar{\mathcal{B}}$ reduces to the following forms for the various strategies (non-cooperative (69a), consensus (82), and diffusion (84a)–(84b)):

$$\bar{\mathcal{B}}_{\text{ncop}} = I_{MN} - \bar{\mathcal{M}}\mathcal{R} \quad (133a)$$

$$\bar{\mathcal{B}}_{\text{cons}} = \bar{\mathcal{A}}^\top - \bar{\mathcal{M}}\mathcal{R} \quad (133b)$$

$$\bar{\mathcal{B}}_{\text{atc}} = \bar{\mathcal{A}}^\top (I_{MN} - \bar{\mathcal{M}}\mathcal{R}) \quad (133c)$$

$$\bar{\mathcal{B}}_{\text{cta}} = (I_{MN} - \bar{\mathcal{M}}\mathcal{R}) \bar{\mathcal{A}}^\top \quad (133d)$$

where $\bar{\mathcal{A}} = \bar{\mathcal{A}} \otimes I_M$ and $\bar{\mathcal{A}} = \mathbb{E} \mathcal{A}_i$. ♦

Example VII.3 (MSE networks with uniform agents). The results of Example VII.2 simplify when all agents employ step-sizes with the same mean value, $\bar{\mu}_k \equiv \bar{\mu}$, and observe regression data with the same covariance matrix, $R_{u,k} \equiv R_u$ [15], [32]. In this case, we can express $\bar{\mathcal{M}}$ and \mathcal{R} from (127b) and (130c) in Kronecker product form as follows:

$$\bar{\mathcal{M}} = \bar{\mu} I_N \otimes I_M, \quad \mathcal{R} = I_N \otimes 2R_u \quad (134)$$

so that expressions (133a)–(133d) reduce to

$$\bar{\mathcal{B}}_{\text{ncop}} = I_N \otimes (I_M - 2\bar{\mu}R_u) \quad (135a)$$

$$\bar{\mathcal{B}}_{\text{cons}} = \bar{\mathcal{A}}^\top \otimes I_M - 2\bar{\mu}(I_M \otimes R_u) \quad (135b)$$

$$\bar{\mathcal{B}}_{\text{atc}} = \bar{\mathcal{A}}^\top \otimes (I_M - 2\bar{\mu}R_u) \quad (135c)$$

$$\bar{\mathcal{B}}_{\text{cta}} = \bar{\mathcal{A}}^\top \otimes (I_M - 2\bar{\mu}R_u) \quad (135d)$$

Observe that $\bar{\mathcal{B}}_{\text{atc}} = \bar{\mathcal{B}}_{\text{cta}}$, so we denote these matrices by $\bar{\mathcal{B}}_{\text{diff}}$. Using properties of the eigenvalues of Kronecker products of matrices, it can be easily verified that the MN eigenvalues of the above $\bar{\mathcal{B}}$ matrices are given by the following expressions in terms of the eigenvalues of the component matrices $\{\bar{\mathcal{A}}, R_u\}$ for $k = 1, 2, \dots, N$ and $m = 1, 2, \dots, M$:

$$\lambda(\bar{\mathcal{B}}_{\text{ncop}}) = 1 - 2\bar{\mu}\lambda_m(R_u) \quad (136a)$$

$$\lambda(\bar{\mathcal{B}}_{\text{cons}}) = \lambda_k(\bar{\mathcal{A}}) - 2\bar{\mu}\lambda_m(R_u) \quad (136b)$$

$$\lambda(\bar{\mathcal{B}}_{\text{diff}}) = \lambda_k(\bar{\mathcal{A}}) [1 - 2\bar{\mu}\lambda_m(R_u)] \quad (136c)$$

♦

Example VII.4 (Potential instability in consensus networks). Consensus strategies can become unstable when used for adaptation purposes [2], [32]. This undesirable effect is already reflected in

expressions (136a)–(136c). In particular, observe that the eigenvalues of $\bar{\mathcal{A}}$ appear multiplying $(1 - 2\bar{\mu}\lambda_m(R_u))$ in expression (136c) for diffusion. As such, and since $\rho(\bar{\mathcal{A}}) = 1$ for any left-stochastic matrix, we conclude for this case of uniform agents that $\rho(\bar{\mathcal{B}}_{\text{diff}}) = \rho(\bar{\mathcal{B}}_{\text{ncop}})$. It follows that, regardless of the choice of the mean combination policy $\bar{\mathcal{A}}$, the diffusion strategies will be stable in the mean (i.e., $\mathbb{E} \tilde{\mathbf{w}}_i$ will converge asymptotically to zero) whenever the individual non-cooperative agents are stable in the mean:

$$\text{individual agents stable} \implies \text{diffusion networks stable} \quad (137a)$$

The same conclusion is not true for consensus networks; the individual agents can be stable and yet the consensus network can become unstable. This is because $\lambda_k(\bar{\mathcal{A}})$ appears as an additive (rather than multiplicative) term in (136b) (see [2], [14], [32] for examples):

$$\text{individual agents stable} \not\Rightarrow \text{consensus networks stable} \quad (137b)$$

The fact that the combination matrix $\bar{\mathcal{A}}^\top$ appears in an additive form in (133b) is the result of the asymmetry that was mentioned earlier in the update equation for the consensus strategy. In contrast, the update equations for the diffusion strategies lead to $\bar{\mathcal{A}}^\top$ appearing in a multiplicative form in (133c)–(133d). ♦

B. Diffusion Networks

Given the superior stability properties of diffusion strategies for adaptation and learning over networks, we continue our presentation by focusing on this class of algorithms. The results in the previous section focus on MSE networks, which deal with mean-square-error cost functions. We now consider networks with more general costs, $\{J_k(w)\}$, and apply diffusion strategies to seek the global minimizer, w^o , of the aggregate cost function, $J^{\text{glob}}(w)$, defined by (57). Without loss in generality, we consider ATC diffusion implementations of the form (99):

$$\psi_{k,i} = \mathbf{w}_{k,i-1} - \boldsymbol{\mu}_k(i) \widehat{\nabla_{w^\top} J_k}(\mathbf{w}_{k,i-1}) \quad (138a)$$

$$\mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_{k,i}} \mathbf{a}_{\ell k}(i) \psi_{\ell,i} \quad (138b)$$

Similar conclusions will apply to CTA diffusion implementations.

Assumption VII.1 (Conditions on cost functions). The aggregate cost $J^{\text{glob}}(w)$ in (57) is twice differentiable and satisfies a condition similar to (9) in Assumption II.1 for some positive parameters $\nu_d \leq \delta_d$. Moreover, all individual costs $\{J_k(w)\}$ are assumed to be strongly-convex with their global minimizers located at w^o , as indicated earlier by (74). □

As explained before following (74), references [1], [2], [49] present results on the case in which the individual costs are only convex and need not be strongly convex. These references also discuss the case in which the individual costs need not share minimizers.

With each agent k in (99), we again associate a gradient noise vector:

$$\mathbf{s}_{k,i}(\mathbf{w}_{k,i-1}) \triangleq \widehat{\nabla_{w^\top} J_k}(\mathbf{w}_{k,i-1}) - \nabla_{w^\top} J_k(\mathbf{w}_{k,i-1}) \quad (139)$$

Assumption VII.2 (Conditions on gradient noise). It is assumed that the first and second-order conditional moments of the gradient

noise components satisfy:

$$\mathbb{E} [\mathbf{s}_{k,i}(\mathbf{w}_{k,i-1}) | \mathcal{F}_{i-1}] = 0 \quad (140a)$$

$$\mathbb{E} [\mathbf{s}_{k,i}(\mathbf{w}_{k,i-1}) \mathbf{s}_{\ell,i}^\top(\mathbf{w}_{\ell,i-1}) | \mathcal{F}_{i-1}] = 0, \quad \forall k \neq \ell \quad (140b)$$

$$\mathbb{E} [\|\mathbf{s}_{k,i}(\mathbf{w}_{k,i-1})\|^2 | \mathcal{F}_{i-1}] \leq \beta_k^2 \|\tilde{\mathbf{w}}_{k,i-1}\|^2 + \sigma_{s,k}^2 \quad (140c)$$

almost surely for some nonnegative scalars β_k^2 and $\sigma_{s,k}^2$, and where \mathcal{F}_{i-1} represents the collection of all random events generated by the iterates from across all agents, $\{\mathbf{w}_{\ell,j}, \ell = 1, 2, \dots, N\}$, up to time $i-1$. Moreover, it is assumed that the limiting covariance matrix of $\mathbf{s}_{k,i}(w^o)$ exists:

$$R_{s,k} \triangleq \lim_{i \rightarrow \infty} \mathbb{E} [\mathbf{s}_{k,i}(w^o) \mathbf{s}_{k,i}^\top(w^o) | \mathcal{F}_{i-1}] \quad (140d)$$

□

We collect the error vectors and gradient noises from across all agents into $N \times 1$ block vectors, whose individual entries are of size $M \times 1$ each:

$$\tilde{\mathbf{w}}_i \triangleq \begin{bmatrix} \tilde{\mathbf{w}}_{1,i} \\ \tilde{\mathbf{w}}_{2,i} \\ \vdots \\ \tilde{\mathbf{w}}_{N,i} \end{bmatrix}, \quad \mathbf{s}_i \triangleq \begin{bmatrix} \mathbf{s}_{1,i} \\ \mathbf{s}_{2,i} \\ \vdots \\ \mathbf{s}_{N,i} \end{bmatrix} \quad (141)$$

and where we are dropping the argument $\mathbf{w}_{k,i-1}$ from the $\mathbf{s}_{k,i}(\cdot)$ for compactness of notation. Likewise, we introduce the following $N \times N$ block diagonal matrices, whose individual entries are of size $M \times M$ each:

$$\mathcal{M}_i = \text{diag}\{\mu_1(i)I_M, \mu_2(i)I_M, \dots, \mu_N(i)I_M\} \quad (142a)$$

$$\mathcal{H}_{i-1} = \text{diag}\{\mathbf{H}_{1,i-1}, \mathbf{H}_{2,i-1}, \dots, \mathbf{H}_{N,i-1}\} \quad (142b)$$

where

$$\mathbf{H}_{k,i-1} \triangleq \int_0^1 \nabla_w^2 J_k(w^o - t\tilde{\mathbf{w}}_{k,i-1}) dt \quad (142c)$$

Now, in a manner similar to (29b), we can appeal to the mean-value theorem [2], [42], [61] to note that

$$\nabla_{w^\top} J_k(\mathbf{w}_{k,i-1}) = -\mathbf{H}_{k,i-1} \tilde{\mathbf{w}}_{k,i-1} \quad (143)$$

so that the approximate gradient vector can be expressed as:

$$\widehat{\nabla_{w^\top} J_k}(\mathbf{w}_{k,i-1}) = -\mathbf{H}_{k,i-1} \tilde{\mathbf{w}}_{k,i-1} + \mathbf{s}_{k,i}(\mathbf{w}_{k,i-1}) \quad (144)$$

Subtracting w^o from both sides of (138a)–(138b), and using (144), we find that the network error vector evolves according to the following stochastic recursion:

$$\tilde{\mathbf{w}}_i = \mathcal{B}_{i-1} \tilde{\mathbf{w}}_{i-1} + \mathcal{A}_i^\top \mathcal{M}_i \mathbf{s}_i \quad (145a)$$

where

$$\mathcal{B}_{i-1} \triangleq \mathcal{A}_i^\top (I_{NM} - \mathcal{M}_i \mathcal{H}_{i-1}) \quad (145b)$$

Recursion (145a) describes the evolution of the network error vector for general convex costs, $J_k(w)$, in a manner similar to recursion (130a) in the mean-square-error case. However, recursion (145a) is more challenging to deal with because of the presence of the random matrix \mathcal{H}_{i-1} ; this matrix is replaced by the constant term \mathcal{R} in the earlier recursion (130a) because that example deals with MSE networks where the

individual costs, $J_k(w)$, are quadratic in w and, therefore, their Hessian matrices are constant and independent of w . In that case, each matrix $\mathbf{H}_{k,i-1}$ in (142c) will evaluate to $2R_{u,k}$ and the matrix \mathcal{H}_{i-1} in (145b) will coincide with the matrix \mathcal{R} defined by (130c).

The next statement ascertains that sufficiently small step-sizes exist that guarantee the MSE stability of the asynchronous diffusion strategy (138a)–(138b) [2], [3].

Lemma VII.1 (MSE network stability). *Consider an asynchronous network of N interacting agents running the ATC diffusion strategy (138a)–(138b). Assume the conditions in Assumptions V.1, VII.1, and VII.2 hold. Let*

$$\mu_{x,\max} = \max_{1 \leq k \leq N} \{\mu_{x,k}\} = \max_{1 \leq k \leq N} \left(\bar{\mu}_k + \frac{\sigma_{\mu,k}^2}{\bar{\mu}_k} \right) \quad (146)$$

Then, there exists $\mu_o > 0$ such that for all $\mu_{x,\max} < \mu_o$:

$$\limsup_{i \rightarrow \infty} \mathbb{E} \|\tilde{\mathbf{w}}_{k,i}\|^2 = O(\mu_{x,\max}) \quad (147)$$

□

Proof: See App. IV of [3].

■

Result (147) shows that the MSD of the network is in the order of $\mu_{x,\max}$. Therefore, sufficiently small step-sizes lead to sufficiently small MSDs. As was the case with the discussion in subsection II-G, we can also seek a closed-form expression for the MSD performance of the asynchronous diffusion network and its agents. To do that, we first introduce the analog of Assumption II.4 for the network case.

Assumption VII.3 (Smoothness Conditions). *The Hessian matrix of the individual cost functions $\{J_k(w)\}$, and the noise covariance matrices defined for each agent in a manner similar to (17a) and denoted by $R_{s,k,i}(w)$, are assumed to be locally Lipschitz continuous in a small neighborhood around w^o :*

$$\|\nabla_w^2 J_k(w^o + \delta w) - \nabla_w^2 J_k(w^o)\| \leq \tau_{k,d} \|\delta w\| \quad (148a)$$

$$\|R_{s,k,i}(w^o + \delta w) - R_{s,k,i}(w^o)\| \leq \tau_{k,s} \|\delta w\|^\kappa \quad (148b)$$

for small perturbation $\|\delta w\| \leq \tau_d$ and for some $\tau_{k,d}, \tau_{k,s} \geq 0$ and $1 \leq \kappa \leq 2$.

Lemma VII.2 (Asynchronous network MSD performance). *Consider an asynchronous network of N interacting agents running the asynchronous diffusion strategy (138a)–(138b). Assume the conditions under Assumptions V.1, VII.1, VII.2, and V.1 hold. Assume further that the step-size parameter $\mu_{x,\max}$ is sufficiently small to ensure mean-square stability, as already ascertained by Lemma VII.1. Then,*

$$\text{MSD}_{\text{diff},k}^{\text{asyn}} \approx \text{MSD}_{\text{diff},\text{av}}^{\text{asyn}} \approx \quad (149a)$$

$$\frac{1}{2} \text{Tr} \left[\left(\sum_{k=1}^N \bar{\mu}_k \bar{p}_k H_k \right)^{-1} \left(\sum_{k=1}^N (\bar{\mu}_k^2 + \sigma_{\mu,k}^2) p_{c,k} R_{s,k} \right) \right]$$

where $H_k = \nabla_w^2 J_k(w^o)$. Moreover, for large enough i , the convergence rate towards the above steady-state value is well approximated by the scalar:

$$\alpha_{\text{dist}}^{\text{asyn}} = 1 - 2\lambda_{\min} \left(\sum_{k=1}^N \bar{\mu}_k \bar{p}_k H_k \right) + O(\mu_{x,\max}^{1+1/N^2}) \quad (149b)$$

□

Proof: See App. XII in [4].

VIII. CONCLUDING REMARKS

This chapter provides an overview of asynchronous strategies for adaptation, learning, and optimization over networks including non-cooperative, centralized, consensus, and diffusion strategies. Particular attention is given to the constant step-size case in order to examine solutions that are able to adapt and learn continuously from streaming data. The presentation complements the results from [1], [2]. We introduced a fairly general model for asynchronous behavior that allows for random step-sizes, link failures, random topology variations, and random combination coefficients. We examined the mean-square-error performance and stability properties under asynchronous events and recovered results for synchronous operation as a special case. The results indicate that asynchronous networks are robust, resilient to failure, and remain mean-square stable for sufficiently small step-sizes.

There are of course several other aspects of distributed strategies that are not covered in this work. Comments on these aspects can be found in [1], [2], [15], including issues related to (a) the noisy exchange of information over links (e.g., [15], [22], [126]–[131]); (b) the use of gossip strategies (e.g., [20], [23], [28], [30], [83], [132], [133]); (c) the exploitation of sparsity constraints (e.g., [134]–[137]); (d) the solution of constrained optimization problems (e.g., [21], [120], [138]–[140]); (e) the use of distributed solutions of the recursive least-squares type (e.g., [15], [84], [118]); (f) the development of distributed state-space solutions (e.g., [84], [104], [109], [117], [141]–[145]); and (g) the study of incremental-based strategies (e.g., [146]–[158]).

REFERENCES

- [1] A. H. Sayed, "Adaptive networks," *Proceedings of the IEEE*, vol. 102, no. 4, pp. 460–497, April 2014.
- [2] A. H. Sayed, *Adaptation, Learning, and Optimization over Networks*, Foundations and Trends in Machine Learning, vol. 7, issue 4–5, pp. 311–801, NOW Publishers, Boston-Delft, July 2014.
- [3] X. Zhao and A. H. Sayed, "Asynchronous adaptation and learning over networks — Part I: Modeling and stability analysis," *IEEE Trans. Signal Processing*, vol. 63, no. 4, pp. 811–826, Feb. 2015.
- [4] X. Zhao and A. H. Sayed, "Asynchronous adaptation and learning over networks — Part II: Performance analysis," *IEEE Trans. Signal Processing*, vol. 63, no. 4, pp. 827–842, Feb. 2015.
- [5] X. Zhao and A. H. Sayed, "Asynchronous adaptation and learning over networks — Part III: Comparison analysis," *IEEE Trans. Signal Processing*, vol. 63, no. 4, pp. 843–858, Feb. 2015.
- [6] S. Camazine, J. L. Deneubourg, N. R. Franks, J. Sneyd, G. Theraulaz, and E. Bonabeau, *Self-Organization in Biological Systems*, Princeton University Press, 2003.
- [7] M. E. J. Newman, "The structure and function of complex networks," *SIAM Review*, vol. 45, no. 2, pp. 167–256, 2003.
- [8] M. E. J. Newman, *Networks: An Introduction*, Oxford University Press, 2010.
- [9] T. G. Lewis, *Network Science: Theory and Applications*, Wiley, NJ, 2009.
- [10] C. W. Reynolds, "Flocks, herds, and schools: A distributed behavior model," *ACM Proc. Comput. Graphs Interactive Tech.*, pp. 25–34, 1987.
- [11] I. D. Couzin, "Collective cognition in animal groups," *Trends in Cognitive Sciences*, vol. 13, pp. 36–43, Jan. 2009.
- [12] O. Sporns, *Networks of the Brain*, MIT Press, 2010.
- [13] S. Haykin, *Cognitive Dynamic Systems*, Cambridge University Press, 2012.
- [14] A. H. Sayed, S.-Y. Tu, J. Chen, X. Zhao, and Z. Towfic, "Diffusion strategies for adaptation and learning over networks," *IEEE Signal Processing Magazine*, vol. 30, no. 3, pp. 155–171, May 2013.
- [15] A. H. Sayed, "Diffusion adaptation over networks," in *E-Reference Signal Processing*, vol. 3, R. Chellapa and S. Theodoridis, Eds., pp. 323–454, Academic Press, 2014. Also available as arXiv:1205.4220v1 [cs.MA], May 2012.
- [16] P. Di Lorenzo, S. Barbarossa, and A. H. Sayed, "Bio-inspired decentralized radio access based on swarming mechanisms over adaptive networks," *IEEE Trans. Signal Processing*, vol. 61, no. 12, pp. 3183–3197, June 2013.
- [17] P. Chainais and C. Richard, "Learning a common dictionary over a sensor network," *Proc. IEEE CAMSAP*, pp. 1–5, Saint Martin, Dec. 2013.
- [18] J. Chen, A. H. Sayed, and Z. Towfic, "Dictionary learning over distributed models," *IEEE Trans. Signal Process.*, vol. 63, issue 4, pp. 1001–1016, February 2015.
- [19] J. Tsitsiklis, D. Bertsekas, and M. Athans, "Distributed asynchronous deterministic and stochastic gradient optimization algorithms," *IEEE Trans. Autom. Control*, vol. 31, no. 9, pp. 803–812, Sep. 1986.
- [20] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah, "Randomized gossip algorithms," *IEEE Trans. on Information Theory*, vol. 52, no. 6, pp. 2508–2530, Jun. 2006.
- [21] K. Srivastava and A. Nedic, "Distributed asynchronous constrained stochastic optimization," *IEEE J. Sel. Topics. Signal Process.*, vol. 5, no. 4, pp. 772–790, Aug. 2011.
- [22] S. Kar and J. M. F. Moura, "Distributed consensus algorithms in sensor networks: Link failures and channel noise," *IEEE Trans. Signal Process.*, vol. 57, no. 1, pp. 355–369, Jan. 2009.
- [23] S. Kar and J. M. F. Moura, "Convergence rate analysis of distributed gossip (linear parameter) estimation: Fundamental limits and tradeoffs," *IEEE Journal on Selected Topics in Signal Processing*, vol. 5, no. 4, pp. 674–690, Aug. 2011.
- [24] S. Kar and J. M. F. Moura, "Sensor networks with random links: Topology design for distributed consensus," *IEEE Trans. Signal Process.*, vol. 56, no. 7, pp. 3315–3326, July 2008.
- [25] D. Jakovetic, J. Xavier, and J. M. F. Moura, "Weight optimization for consensus algorithms with correlated switching topology," *IEEE Trans. Signal Process.*, vol. 58, no. 7, pp. 3788–3801, July 2010.
- [26] D. Jakovetic, J. Xavier, and J. M. F. Moura, "Cooperative convex optimization in networked systems: Augmented Lagrangian algorithms with directed Gossip communication," *IEEE Trans. Signal Process.*, vol. 59, no. 8, pp. 3889–3902, Aug. 2011.
- [27] S. Kar and J. M. F. Moura, "Distributed consensus algorithms in sensor networks: Quantized data and random link failures," *IEEE Trans. Signal Process.*, vol. 58, no. 3, pp. 1383–1400, Mar. 2010.
- [28] T. C. Aysal, M. E. Yildiz, A. D. Sarwate, and A. Scaglione, "Broadcast gossip algorithms for consensus," *IEEE Trans. on Signal Processing*, vol. 57, no. 7, pp. 2748–2761, July 2009.
- [29] T. C. Aysal, A. D. Sarwate, and A. G. Dimakis, "Reaching consensus in wireless networks with probabilistic broadcast," in *Proc. Allerton Conf. Commun., Control, Comput.*, Allerton House, IL, Sept. and Oct. 2009, pp. 732–739.
- [30] C. Lopes and A. H. Sayed, "Diffusion adaptive networks with changing topologies," *Proc. IEEE ICASSP*, pp. 3285–3288, Las Vegas, April 2008.
- [31] N. Takahashi and I. Yamada, "Link probability control for probabilistic diffusion least-mean squares over resource-constrained networks," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, pp. 3518–3521, Dallas, TX, Mar. 2010.
- [32] S.-Y. Tu and A. H. Sayed, "Diffusion strategies outperform consensus strategies for distributed estimation over adaptive networks," *IEEE Trans. on Signal Processing*, vol. 60, no. 12, pp. 6217–6234, Dec. 2012.
- [33] S. Haykin, *Adaptive Filter Theory*, Prentice Hall, NJ, 2002.
- [34] B. Widrow and S. D. Stearns, *Adaptive Signal Processing*, Prentice Hall, NJ, 1985.
- [35] A. H. Sayed, *Adaptive Filters*, Wiley, NJ, 2008.
- [36] A. H. Sayed, *Fundamentals of Adaptive Filtering*, Wiley, NJ, 2003.
- [37] T. Kailath, A. H. Sayed, and B. Hassibi, *Linear Estimation*, Prentice Hall, NJ, 2000.
- [38] C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2007.
- [39] S. Theodoridis and K. Koutroumbas, *Pattern Recognition*, 4th edition, Academic Press, 2008.
- [40] D. W. Hosmer and S. Lemeshow, *Applied Logistic Regression*, 2nd edition, Wiley, NJ, 2000.
- [41] B. T. Poljak and Y. Z. Tsytkin, "Pseudogradient adaptation and training algorithms," *Autom. Remote Control*, vol. 12, pp. 83–94, 1973.

- [42] B. Poljak, *Introduction to Optimization*, Optimization Software, NY, 1987.
- [43] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*, 1st edition, Athena Scientific, Singapore, 1997.
- [44] Y. Z. Tsypkin, *Adaptation and Learning in Automatic Systems*, Academic Press, NY, 1971.
- [45] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
- [46] D. Bertsekas, *Convex Analysis and Optimization*, Athena Scientific, 2003.
- [47] Y. Nesterov, *Introductory Lectures on Convex Optimization: A Basic Course*, Kluwer Academic Publishers, 2004.
- [48] J. Chen and A. H. Sayed, "Diffusion adaptation strategies for distributed optimization and learning over networks," *IEEE Trans. Signal Processing*, vol. 60, no. 8, pp. 4289–4305, Aug. 2012.
- [49] J. Chen and A. H. Sayed, "On the learning behavior of adaptive networks — Part I: Transient analysis," *submitted for publication*. Also available as arXiv:1312.7581 [cs.MA], Dec. 2013.
- [50] H. Robbins and S. Monro, "A stochastic approximation method," *Ann. Math. Stat.*, vol. 22, pp. 400–407, 1951.
- [51] R. von Mises and H. Pollaczek-Geiringer, "Praktische verfahren der gleichungs-auflösung," *Z. Angew. Math. Mech.*, vol. 9, pp. 152–164, 1929.
- [52] J. R. Blum, "Multidimensional stochastic approximation methods," *Ann. Math. Stat.*, vol. 25, pp. 737–744, 1954.
- [53] L. Schmetterer, "Stochastic approximation," *Proc. Berkeley Symp. Math. Statist. Probab.*, pp. 587–609, 1961.
- [54] G. B. Wetherhill, *Sequential Methods in Statistics*, Methuen, London, 1966.
- [55] B. Widrow and M. E. Hoff, Jr., "Adaptive switching circuits," *IRE WESCON Conv. Rec.*, Pt. 4, pp. 96–104, 1960.
- [56] D. P. Bertsekas and J. N. Tsitsiklis, "Gradient convergence in gradient methods with errors," *SIAM J. Optim.*, vol. 10, no. 3, pp. 627–642, 2000.
- [57] W. Feller, *An Introduction to Probability Theory and Its Applications*, vol. 2, Wiley, NY, 1971.
- [58] G. J. Hahn and S. Shapiro, *Statistical Models in Engineering*, Wiley, NJ, 1994.
- [59] M. Abramowitz and I. Stegun, Eds., *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, Dover, NY, 1972.
- [60] G. E. Andrews, R. Askey and R. Roy, *Special Functions*, Cambridge University Press, Cambridge, 1999.
- [61] W. Rudin, *Principles of Mathematical Analysis*, McGraw-Hill, 1976.
- [62] N. R. Yousef and A. H. Sayed, "A unified approach to the steady-state and tracking analysis of adaptive filters," *IEEE Trans. Signal Processing*, vol. 49, no. 2, pp. 314–324, February 2001.
- [63] T. Y. Al-Naffouri and A. H. Sayed, "Transient analysis of data-normalized adaptive filters," *IEEE Trans. Signal Processing*, vol. 51, no. 3, pp. 639–652, Mar. 2003.
- [64] J. Chen and A. H. Sayed, "On the learning behavior of adaptive networks — Part II: Performance analysis," *submitted for publication*. Also available as arXiv:1312.7580 [cs.MA], Dec. 2013.
- [65] A. Papoulis and S. U. Pillai, *Probability, Random Variables, and Stochastic Processes*, McGraw-Hill, NY, 2002.
- [66] R. Durrett, *Probability Theory and Examples*, 2nd edition, Duxbury Press, 1996.
- [67] R. M. Dudley, *Real Analysis and Probability*, 2nd edition, Cambridge University Press, 2003.
- [68] B. Widrow, J. M. McCool, M. G. Larimore, and C. R. Johnson Jr., "Stationary and nonstationary learning characteristics of the LMS adaptive filter," *Proc. IEEE*, vol. 64, no. 8, pp. 1151–1162, Aug. 1976.
- [69] L. Horowitz and K. Senne, "Performance advantage of complex LMS for controlling narrow-band adaptive arrays," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 29, no. 3, pp. 722–736, Jun. 1981.
- [70] S. Jones, R. C. III, and W. Reed, "Analysis of error-gradient adaptive linear estimators for a class of stationary dependent processes," *IEEE Trans. Inf. Theory*, vol. 28, no. 2, pp. 318–329, Mar. 1982.
- [71] W. A. Gardner, "Learning characteristics of stochastic-gradient-descent algorithms: A general study, analysis, and critique," *Signal Process.*, vol. 6, no. 2, pp. 113–133, Apr. 1984.
- [72] A. Feuer and E. Weinstein, "Convergence analysis of LMS filters with uncorrelated Gaussian data," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 33, no. 1, pp. 222–230, Feb. 1985.
- [73] J. B. Foley and F. M. Boland, "A note on the convergence analysis of LMS adaptive filters with Gaussian data," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 36, no. 7, pp. 1087–1089, Jul. 1988.
- [74] J. Chen and A. H. Sayed, "Distributed Pareto optimization via diffusion strategies," *IEEE J. Selected Topics in Signal Processing*, vol. 7, no. 2, pp. 205–220, April 2013.
- [75] V. N. Vapnik, *The Nature of Statistical Learning Theory*, Springer, NY, 2000.
- [76] Z. Towfic, J. Chen, and A. H. Sayed, "On the generalization ability of distributed online learners," *Proc. IEEE Workshop on Machine Learning for Signal Processing (MLSP)*, Santander, Spain, pp. 1–6, Sep. 2012.
- [77] X. Zhao and A. H. Sayed, "Performance limits for distributed estimation over LMS adaptive networks," *IEEE Trans. Signal Processing*, vol. 60, no. 10, pp. 5107–5124, Oct. 2012.
- [78] L. Xiao and S. Boyd, "Fast linear iterations for distributed averaging," *Syst. Control Lett.*, vol. 53, no. 1, pp. 65–78, Sep. 2004.
- [79] J. Tsitsiklis and M. Athans, "Convergence and asymptotic agreement in distributed decision problems," *IEEE Trans. Autom. Control*, vol. 29, no. 1, pp. 42–50, Jan. 1984.
- [80] A. Nedic and A. Ozdaglar, "Cooperative distributed multi-agent optimization," in *Convex Optimization in Signal Processing and Communications*, Y. Eldar and D. Palomar (Eds.), Cambridge University Press, pp. 340–386, 2010.
- [81] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Trans. Autom. Control*, vol. 54, no. 1, pp. 48–61, Jan. 2009.
- [82] B. Johansson, T. Keviczky, M. Johansson, and K. Johansson, "Subgradient methods and consensus algorithms for solving convex optimization problems," *Proc. IEEE CDC*, pp. 4185–4190, Cancun, Mexico, December 2008.
- [83] A. G. Dimakis, S. Kar, J. M. F. Moura, M. G. Rabbat, and A. Scaglione, "Gossip algorithms for distributed signal processing," *Proc. IEEE*, vol. 98, no. 11, pp. 1847–1864, Nov. 2010.
- [84] L. Xiao, S. Boyd and S. Lall, "A space-time diffusion scheme peer-to-peer least-squares-estimation," *Proc. Information Processing in Sensor Networks (IPSN)*, pp. 168–176, Nashville, TN, April 2006.
- [85] W. Ren and R. W. Beard, "Consensus seeking in multi-agent systems under dynamically changing interaction topologies," *IEEE Trans. on Automatic Control*, vol. 50, pp. 655–661, May 2005.
- [86] R. Olfati-Saber and J. Shamma, "Consensus filters for sensor networks and distributed sensor fusion," *Proc. 44th IEEE Conference on Decision and Control (CDC)*, pp. 6698–6703, Seville, Spain, Dec. 2005.
- [87] S. Barbarossa, and G. Scutari, "Bio-inspired sensor network design," *IEEE Signal Processing Magazine*, vol. 24, no. 3, pp. 26–35, May 2007.
- [88] S. Sardellitti, M. Giona, and S. Barbarossa, "Fast distributed average consensus algorithms based on advection-diffusion processes," *IEEE Trans. Signal Processing*, vol. 58, no. 2, pp. 826–842, Feb. 2010.
- [89] P. Braca, S. Marano, and V. Matta, "Running consensus in wireless sensor networks," *Proc. 11th International Conference on Information Fusion*, pp. 1–6, Cologne, Germany, June 2008.
- [90] C. G. Lopes and A. H. Sayed, "Diffusion least-mean squares over adaptive networks: Formulation and performance analysis," *IEEE Trans. Signal Process.*, vol. 56, no. 7, pp. 3122–3136, July 2008.
- [91] F. S. Cattivelli and A. H. Sayed, "Diffusion LMS strategies for distributed estimation," *IEEE Trans. Signal Process.*, vol. 58, no. 3, pp. 1035–1048, Mar. 2010.
- [92] R. A. Horn and C. R. Johnson, *Matrix Analysis*, Cambridge University Press, 2003.
- [93] S.-Y. Tu and A. H. Sayed, "Mobile adaptive networks," *IEEE J. Sel. Topics. Signal Process.*, vol. 5, no. 4, pp. 649–664, Aug. 2011.
- [94] F. Cattivelli and A. H. Sayed, "Modeling bird flight formations using diffusion adaptation," *IEEE Transactions on Signal Processing*, vol. 59, no. 5, pp. 2038–2051, May 2011.
- [95] O. Dekel, R. Gilad-Bachrach, O. Shamir, and L. Xiao, "Optimal distributed online prediction," *Proc. International Conference on Machine Learning (ICML)*, Bellevue, WA, pp. 713–720, Jun. 2011.
- [96] A. Agarwal and J. Duchi, "Distributed delayed stochastic optimization," *Proc. Neural Information Processing Systems (NIPS)*, Granada, Spain, pp. 873–881, Dec. 2011.
- [97] J. B. Predd, S. B. Kulkarni, and H. V. Poor, "Distributed learning in wireless sensor networks," *IEEE Signal Processing Magazine*, vol. 23, no. 4, pp. 56–69, Jul. 2006.
- [98] Z. J. Towfic, J. Chen, and A. H. Sayed, "Collaborative learning of mixture models using diffusion adaptation," *Proc. IEEE Workshop Mach. Learn. Signal Process. (MLSP)*, Beijing, China, pp. 1–6, Sep. 2011.
- [99] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 3rd edition, The John Hopkins University Press, Baltimore, 1996.
- [100] A. Berman and R. J. Plemmons, *Nonnegative Matrices in the Mathematical Sciences*, SIAM, PA, 1994.

- [101] S. U. Pillai, T. Suel, and S. Cha, "The Perron–Frobenius theorem: Some of its applications," *IEEE Signal Process. Mag.*, vol. 22, no. 2, pp. 62–75, Mar. 2005.
- [102] M. H. DeGroot, "Reaching a consensus," *Journal of the American Statistical Association*, vol. 69, no. 345, pp. 118–121, 1974.
- [103] R. L. Berger, "A necessary and sufficient condition for reaching a consensus using DeGroot's method," *Journal of the American Statistical Association*, vol. 76, no. 374, pp. 415–418, Jun. 1981.
- [104] P. Aliksson and A. Rantzer, "Distributed Kalman filtering using weighted averaging," *Proc. 17th Int. Symp. Math. Thy Net. Sys (MTNS)*, pp. 1–6, Kyoto, Japan, 2006.
- [105] S. Kar, J. M. F. Moura, and K. Ramanan, "Distributed parameter estimation in sensor networks: Nonlinear observation models and imperfect communication," *IEEE Trans. Information Theory*, vol. 58, no. 6, pp. 3575–3605, Jun. 2012.
- [106] R. Olfati-Saber, "Distributed Kalman filter with embedded consensus filters," *Proc. Joint IEEE Conf. on Decision and Control and European Control Conf. (CDC-ECC)*, pp. 8179–8184, Seville, Spain, Dec. 2005.
- [107] A. Das and M. Mesbahi, "Distributed linear parameter estimation in sensor networks based on Laplacian dynamics consensus algorithm," *Proc. IEEE SECON*, vol. 2, pp. 440–449, Reston, VA, Sep. 2006.
- [108] R. Carli, A. Chiuso, L. Schenato, and S. Zampieri, "Distributed Kalman filtering using consensus strategies," *IEEE J. Sel. Areas Communications*, vol. 26, no. 4, pp. 622–633, Sep. 2008.
- [109] U. A. Khan and J. M. F. Moura, "Distributing the Kalman filter for large-scale systems," *IEEE Trans. Signal Processing*, vol. 56, no. 10, pp. 4919–4935, Oct. 2008.
- [110] A. Speranzon, C. Fischione, and K. H. Johansson, "Distributed and collaborative estimation over wireless sensor networks," *Proc. IEEE CDC*, pp. 1025–1030, San Diego, USA, Dec. 2006.
- [111] C. G. Lopes and A. H. Sayed, "Distributed processing over adaptive networks," in *Proc. Adaptive Sensor Array Processing Workshop*, MIT Lincoln Laboratory, MA, pp. 1–5, June 2006.
- [112] A. H. Sayed and C. G. Lopes, "Adaptive processing over distributed networks," *IEICE Trans. Fund. of Electron., Commun. and Comput. Sci.*, vol. E90-A, no. 8, pp. 1504–1510, 2007.
- [113] C. G. Lopes and A. H. Sayed, "Diffusion least-mean-squares over adaptive networks," *Proc. IEEE ICASSP*, Honolulu, Hawaii, vol. 3, pp. 917–920, April 2007.
- [114] C. G. Lopes and A. H. Sayed, "Steady-state performance of adaptive diffusion least-mean squares," *Proc. IEEE Workshop on Statistical Signal Processing (SSP)*, pp. 136–140, Madison, WI, Aug. 2007.
- [115] F. S. Cattivelli, C. G. Lopes, and A. H. Sayed, "A diffusion RLS scheme for distributed estimation over adaptive networks," *Proc. IEEE Workshop on Signal Process. Advances Wireless Comm. (SPAWC)*, Helsinki, Finland, pp. 1–5, June 2007.
- [116] F. S. Cattivelli and A. H. Sayed, "Diffusion LMS algorithms with information exchange," *Proc. Asilomar Conf. Signals, Syst. Comput.*, Pacific Grove, CA, pp. 251–255, Nov. 2008.
- [117] F. S. Cattivelli and A. H. Sayed, "Diffusion mechanisms for fixed-point distributed Kalman smoothing," *Proc. EUSIPCO*, Lausanne, Switzerland, pp. 1–4, Aug. 2008.
- [118] F. S. Cattivelli, C. G. Lopes, and A. H. Sayed, "Diffusion recursive least-squares for distributed estimation over adaptive networks," *IEEE Trans. Signal Process.*, vol. 56, no. 5, pp. 1865–1877, May 2008.
- [119] S. S. Ram, A. Nedic, and V. V. Veeravalli, "Distributed stochastic subgradient projection algorithms for convex optimization," *J. Optim. Theory Appl.*, vol. 147, no. 3, pp. 516–545, 2010.
- [120] S. Lee and A. Nedic, "Distributed random projection algorithm for convex optimization," *IEEE J. Selected Topics in Signal Processing*, vol. 7, no. 2, pp. 221–229, Apr. 2013.
- [121] P. Bianchi, G. Fort, and W. Hachem, "Performance of a distributed stochastic approximation algorithm," *IEEE Trans. Information Theory*, vol. 59, no. 11, pp. 7405–7418, Nov. 2013.
- [122] S. S. Stankovic, M. S. Stankovic, and D. S. Stipanovic, "Decentralized parameter estimation by consensus based stochastic approximation," *IEEE Trans. on Autom. Control*, vol. 56, no. 3, pp. 531–543, Mar. 2011.
- [123] J. Chen and A. H. Sayed, "On the limiting behavior of distributed optimization strategies," *Proc. 50th Annual Allerton Conference on Communication, Control, and Computing*, pp. 1535–1542, Monticello, IL, Oct. 2012.
- [124] J. A. Bondy and U. S. R. Murty, *Graph Theory*, Springer, 2008.
- [125] J. Leskovec, D. Chakrabarti, J. Kleinberg, C. Faloutsos, and Z. Ghahramani, "Kronecker graphs: An approach to modeling networks," vol. 11, pp. 985–1042, Sept. 2010.
- [126] S.-Y. Tu and A. H. Sayed, "Adaptive networks with noisy links," *Proc. IEEE Globecom*, pp. 1–5, Houston, TX, December 2011.
- [127] X. Zhao, S.-Y. Tu, and A. H. Sayed, "Diffusion adaptation over networks under imperfect information exchange and non-stationary data," *IEEE Trans. Signal Processing*, vol. 60, no. 7, pp. 3460–3475, July 2012.
- [128] R. Abdoole and B. Champagne, "Diffusion LMS algorithms for sensor networks over non-ideal inter-sensor wireless channels," *Proc. IEEE Int. Conf. Dist. Comput. Sensor Systems (DCOSS)*, pp. 1–6, Barcelona, Spain, June 2011.
- [129] A. Khalili, M. A. Tinati, A. Rastegarnia, and J. A. Chambers, "Steady state analysis of diffusion LMS adaptive networks with noisy links," *IEEE Trans. Signal Processing*, vol. 60, no. 2, pp. 974–979, Feb. 2012.
- [130] X. Zhao and A. H. Sayed, "Combination weights for diffusion strategies with imperfect information exchange," *Proc. IEEE ICC*, pp. 648–652, Ottawa, Canada, June 2012.
- [131] G. Mateos, I. D. Schizas, and G. B. Giannakis, "Performance analysis of the consensus-based distributed LMS algorithm," *EURASIP J. Adv. Signal Process.*, pp. 1–19, 2009, 10.1155/2009/981030, Article ID 981030.
- [132] D. Shah, "Gossip algorithms," *Found. Trends Netw.*, vol. 3, pp. 1–125, 2009.
- [133] O. L. Rortveit, J. H. Husoy, and A. H. Sayed, "Diffusion LMS with communications constraints," *Proc. 44th Asilomar Conference on Signals, Systems and Computers*, Pacific Grove, CA, pp. 1645–1649, Nov. 2010.
- [134] P. Di Lorenzo and A. H. Sayed, "Sparse distributed learning based on diffusion adaptation," *IEEE Trans. Signal Processing*, vol. 61, no. 6, pp. 1419–1433, March 2013.
- [135] S. Chouvardas, K. Slavakis, Y. Kopsinis, S. Theodoridis, "A sparsity-promoting adaptive algorithm for distributed learning," *IEEE Transactions on Signal Processing*, vol. 60, no. 10, pp. 5412–5425, Oct. 2012.
- [136] S. Chouvardas, G. Mileounis, N. Kalouptsidis, and S. Theodoridis, "A greedy sparsity-promoting LMS for distributed adaptive learning in diffusion networks," *Proc. ICASSP*, pp. 5415–5419, Vancouver, BC, Canada, 2013.
- [137] Y. Liu, C. Li and Z. Zhang, "Diffusion sparse least-mean squares over networks," *IEEE Transactions on Signal Processing*, vol. 60, no. 8, pp. 4480–4485, Aug. 2012.
- [138] F. Yan, S. Sundaram, S. V. N. Vishwanathan, and Y. Qi, "Distributed autonomous online learning: Regrets and intrinsic privacy-preserving properties," *IEEE Trans. Knowledge and Data Engineering*, vol. 25, no. 11, pp. 2483–2493, Nov. 2013.
- [139] S. Theodoridis, K. Slavakis, and I. Yamada, "Adaptive learning in a world of projections: A unifying framework for linear and nonlinear classification and regression tasks," *IEEE Signal Processing Magazine*, vol. 28, no. 1, pp. 97–123, Jan. 2011.
- [140] Z. Towfic and A. H. Sayed, "Adaptive penalty-based distributed stochastic convex optimization," *IEEE Trans. Signal Process.*, vol. 62, no. 15, pp. 3924–3938, August 2014.
- [141] R. Olfati-Saber, "Kalman-consensus filter: Optimality, stability, and performance," *Proc. IEEE CDC*, pp. 7036–7042, Shanghai, China, 2009.
- [142] R. Olfati-Saber, "Distributed Kalman filtering for sensor networks," *Proc. 46th IEEE Conf. Decision Control*, pp. 5492–5498, New Orleans, LA, Dec. 2007.
- [143] F. S. Cattivelli, C. G. Lopes, and A. H. Sayed, "Diffusion strategies for distributed Kalman filtering: Formulation and performance analysis," *Proc. IAPR Workshop on Cognitive Inf. Process. (CIP)*, Santorini, Greece, pp. 36–41, June 2008.
- [144] F. Cattivelli and A. H. Sayed, "Diffusion strategies for distributed Kalman filtering and smoothing," *IEEE Trans. Automatic Control*, vol. 55, no. 9, pp. 2069–2084, Sep. 2010.
- [145] F. Cattivelli and A. H. Sayed, "Diffusion distributed Kalman filtering with adaptive weights," *Proc. Asilomar Conference on Signals, Systems and Computers*, pp. 908–912, Pacific Grove, CA, Nov. 2009.
- [146] D. P. Bertsekas, "A new class of incremental gradient methods for least squares problems," *SIAM J. Optim.*, vol. 7, no. 4, pp. 913–926, 1997.
- [147] D. P. Bertsekas, *Nonlinear Programming*, 2nd edition, Athena Scientific, Belmont, MA, 1999.
- [148] A. Nedic and D. P. Bertsekas, "Incremental subgradient methods for nondifferentiable optimization," *SIAM J. Optim.*, vol. 12, no. 1, pp. 109–138, 2001.
- [149] M. G. Rabbat and R. D. Nowak, "Quantized incremental algorithms for distributed optimization," *IEEE J. Sel. Areas Commun.*, vol. 23, no. 4, pp. 798–808, 2005.
- [150] E. S. Helou and A. R. De Pierro, "Incremental subgradients for constrained convex optimization: A unified framework and new methods," *SIAM J. on Optimization*, vol. 20, pp. 1547–1572, 2009.

- [151] B. Johansson, M. Rabi, and M. Johansson, "A randomized incremental subgradient method for distributed optimization in networked systems," *SIAM J. on Optimization*, vol. 20, pp. 1157–1170, 2009.
- [152] D. Blatt, A. O. Hero, and H. Gauchman, "A convergent incremental gradient method with a constant step size," *SIAM J. Optimization*, vol. 18, pp. 29–51, 2008.
- [153] A. H. Sayed and C. Lopes, "Distributed recursive least-squares strategies over adaptive networks," *Proc. 40th Asilomar Conference on Signals, Systems and Computers*, Pacific Grove, CA, pp. 233–237, Oct.-Nov., 2006.
- [154] C. G. Lopes and A. H. Sayed, "Incremental adaptive strategies over distributed networks," *IEEE Trans. Signal Process.*, vol. 55, no. 8, pp. 4064–4077, Aug. 2007.
- [155] A. H. Sayed and F. Cattivelli, "Distributed adaptive learning mechanisms," *Handbook on Array Processing and Sensor Networks*, S. Haykin and K. J. Ray Liu, Eds., pp. 695–722, Wiley, NJ, 2009.
- [156] L. Li, J. Chambers, C. G. Lopes, and A. H. Sayed, "Distributed estimation over an adaptive incremental network based on the affine projection algorithm," *IEEE Trans. Signal Processing*, vol. 58, no. 1, pp. 151–164, Jan. 2010.
- [157] F. Cattivelli and A. H. Sayed, "Analysis of spatial and incremental LMS processing for distributed estimation," *IEEE Trans. Signal Processing*, vol. 59, no. 4, pp. 1465–1480, April 2011.
- [158] J. B. Predd, S. R. Kulkarni, and H. V. Poor, "A collaborative training algorithm for distributed learning," *IEEE Trans. Information Theory*, vol. 55, no. 4, pp. 1856–1871, April 2009.